# An Introduction to Interdomain Routing and the Border Gateway Protocol (BGP)

## Timothy G. Griffin

**AT&T Research**

**griffin@research.att.com**

**http://www.research.att.com/~griffin**

**http://www.research.att.com/~griffin/interdomain.html**

## ICNP
## PARIS

**November 12, 2002**

# Outline

1. The Internet is implemented with a diverse set of physical networks
2. Relationships between Autonomous Routing Domains (ARDs)
3. BGP as a means of implementing and maintaining relationships between ARDs
4. BGP as means of implementing local optimizations ("Traffic Engineering")
5. What Problem is BGP Solving anyway?
6. Current Internet Growth Trends
7. Selected References
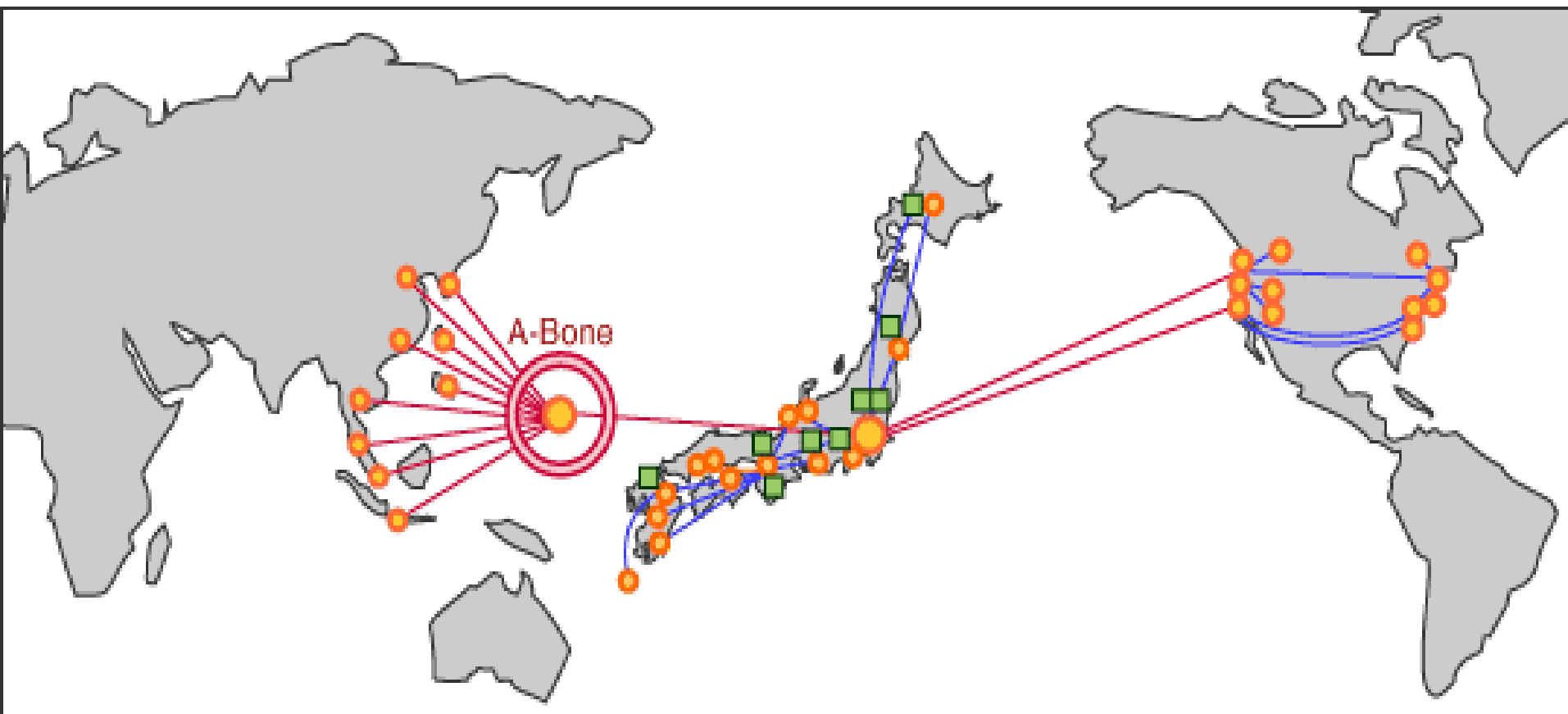
# In Memory of Abha Ahuja



Photo by Peter Lothberg. http://www.caida.org/~kc/abha/gallery.html

NANOG memorial site: http://www.nanog.org/abha.html
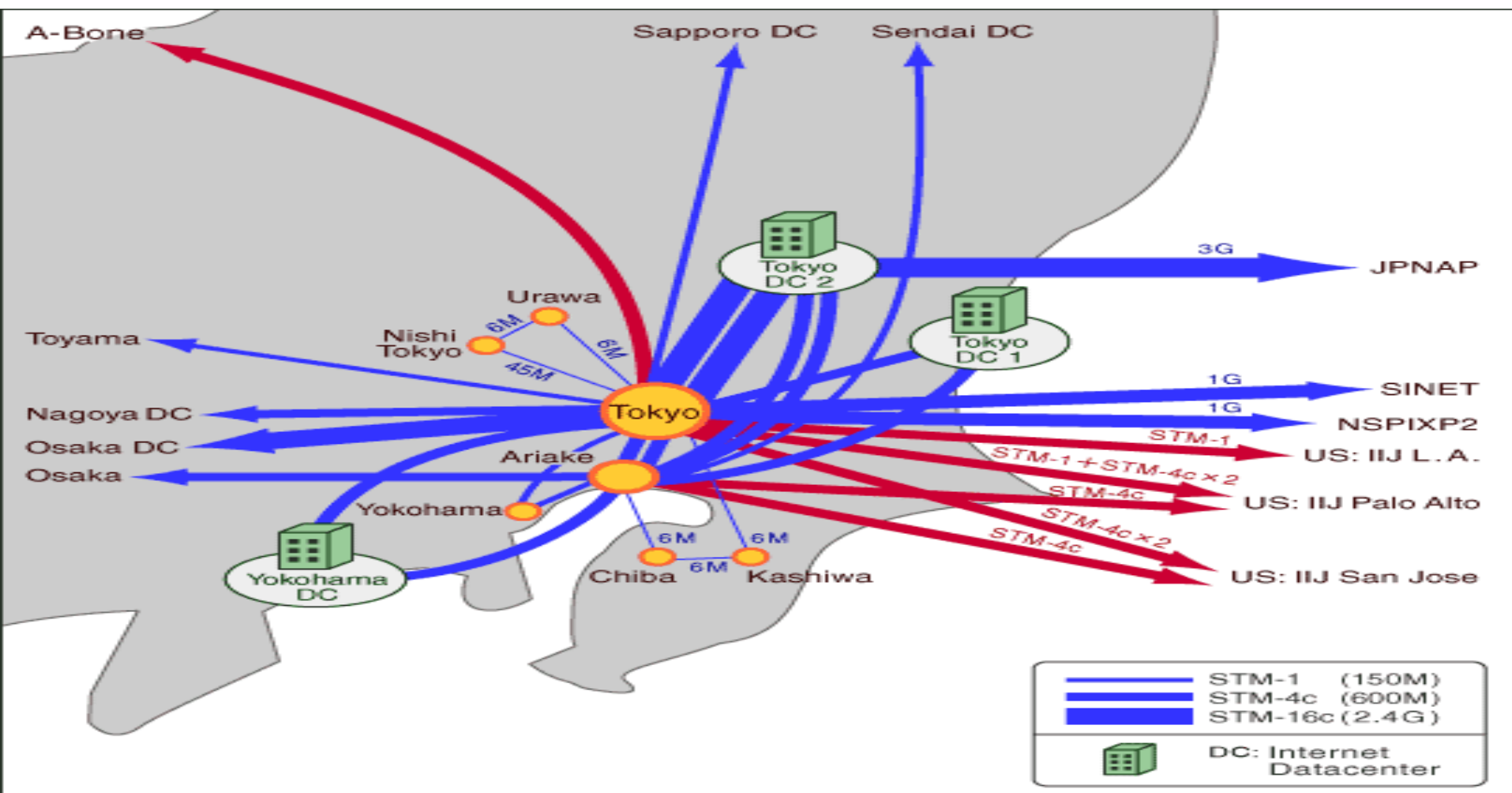
ICNP 2002

# PART I

# Physical Connectivity

# Internet Initiative Japan (IIJ)



A-Bone

# IIJ, Tokyo

# Telstra international



Peering Links
| | |
|---|---|
| Japan | 55 Mbps |
| South Korea | 8 Mbps |
| China | 8 Mbps |
| Taiwan | 4 Mbps |
| Hong Kong | 45 Mbps |
| Malaysia | 2 Mbps |
| Singapore | 45 Mbps |

USA
620 Mbps

USA
445 Mbps

USA
200 Mbps

New Zealand
200 Mbps

# WorldCom (UUNet)



| | |
|---|---|
| —— 64 Kbps | ▬ OC12c/STM4 (622 Mbps) |
| —— T1/E1 (1.5 Mbps/2 Mbps) | ▬ OC48c/STM16 (2.5 Gbps) |
| —— E3/T3/DS3 (35 Mbps/45 Mbps) | ▬ OC192c/STM64 (10 Gbps) |
| —— T2 (6 Mbps) | • Single Hub City |
| —— OC3c/STM1 (155 Mbps) | ■ Multiple Hubs City |
| | ◉ Data Center Hub |

# UUNet, Europe



EMEA detail, JAN 2002

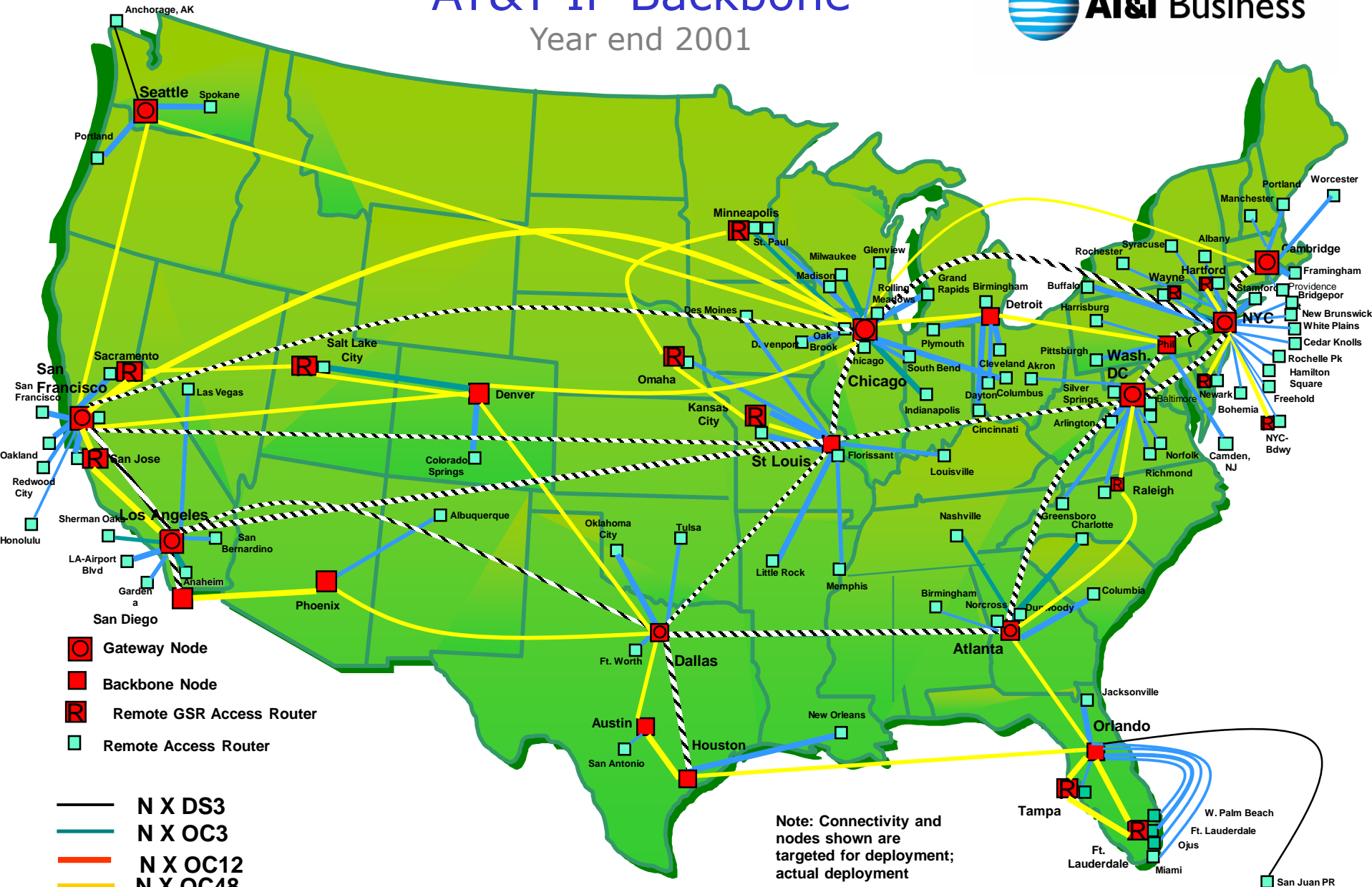| | |
|---|---|
| ——— 64 Kbps | ——— OC12c/STM4 (622 Mbps) |
| ——— T1/E1 (1.5 Mbps/2 Mbps) | ——— OC48c/STM16 (2.5 Gbps) |
| ——— E3/T3/DS3 (35 Mbps/45 Mbps) | ——— OC192c/STM64 (10 Gbps) |
| ——— T2 (6 Mbps) | ● Single Hub City |
| ——— OC3c/STM1 (155 Mbps) | ■ Multiple Hubs City |
| | ◉ Data Center Hub |

# Sprint, USA



U.S. Sprint IP Backbone Network and Internet Centers (Q3 2001)

# AT&T IP Backbone
## Year end 2001

**AT&T Business**

**Legend:**
- ⬤ **Gateway Node**
- ⬛ **Backbone Node**
- **R** **Remote GSR Access Router**
- ⬛ **Remote Access Router**

**Link types:**
- — **N X DS3**
- — **N X OC3**
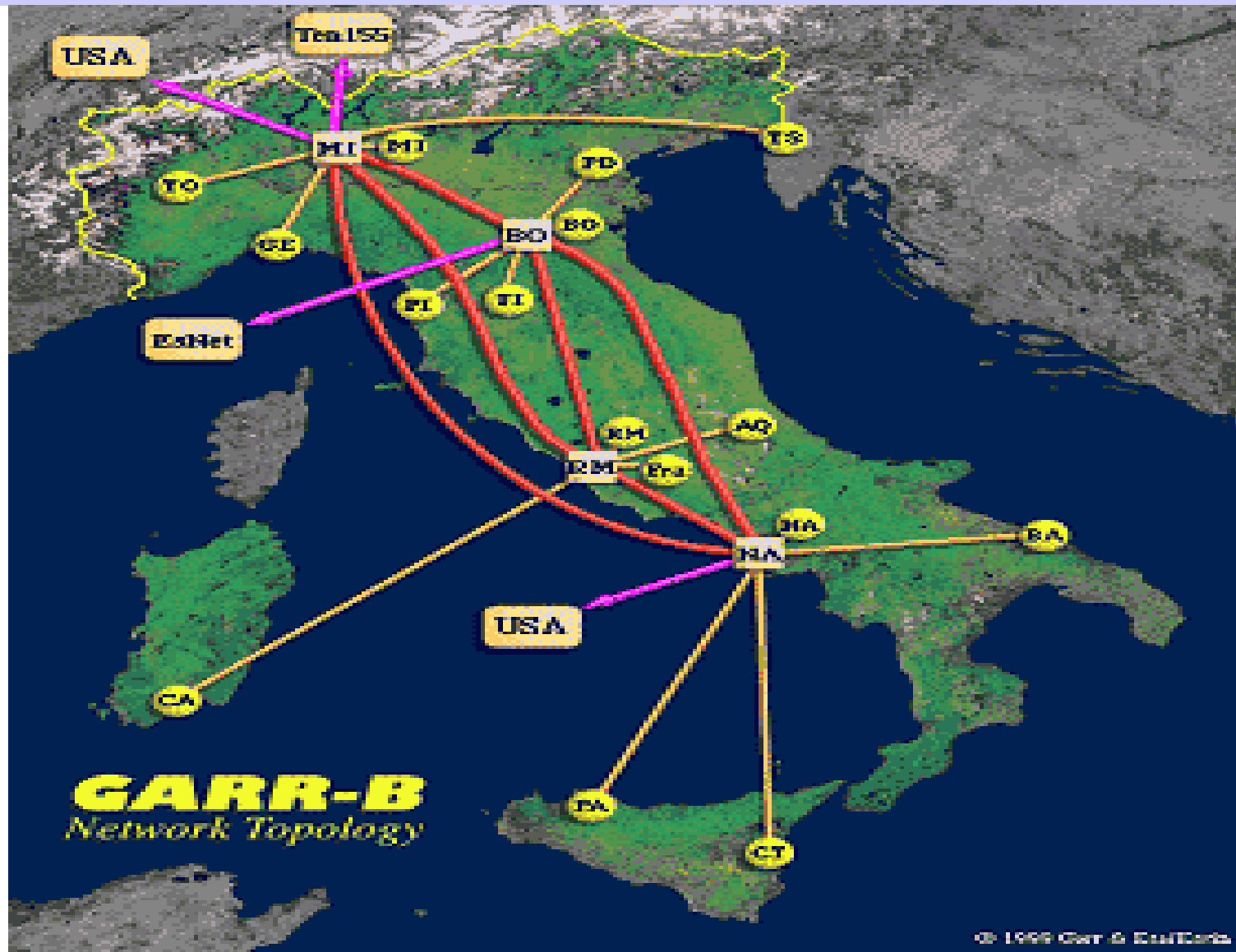- — **N X OC12**
- — **N X OC48**
- ▨ **NX OC192**

Note: Connectivity and nodes shown are targeted for deployment; actual deployment may vary. Maps should not be used to predict service availability.
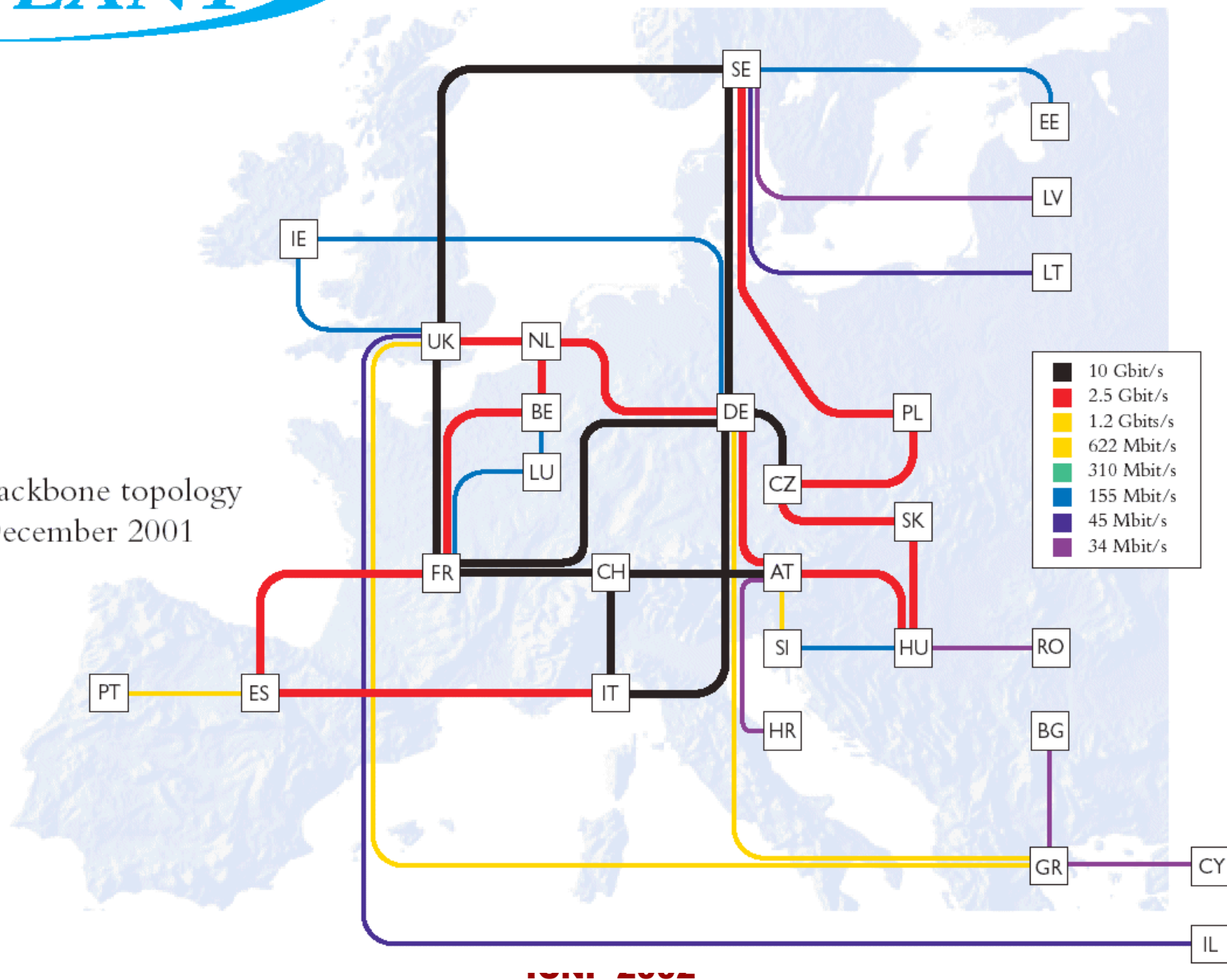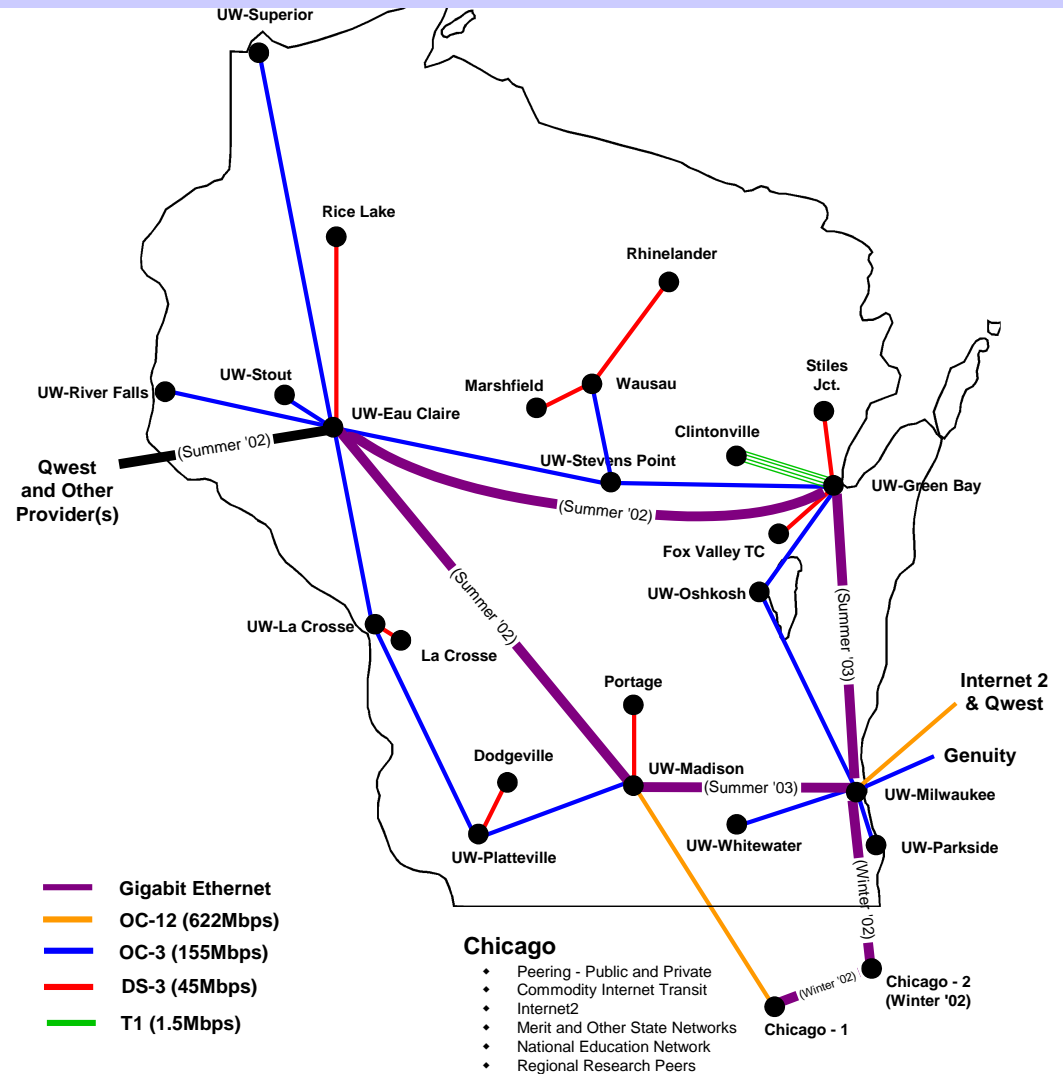
**ICNP 2002**

Rev. 6-4-01

**City labels:**
Anchorage, AK · Seattle · Spokane · Portland · Minneapolis · St. Paul · Milwaukee · Glenview · Madison · Rolling Meadows · Grand Rapids · Birmingham · Buffalo · Rochester · Syracuse · Albany · Cambridge · Framingham · Manchester · Portland · Worcester · Wayne · Hartford · Stamford · Providence · Bridgepor · New Brunswick · White Plains · Cedar Knolls · Rochelle Pk · Hamilton Square · Freehold · NYC · Des Moines · Davenport · Oak Brook · Chicago · Plymouth · Detroit · Harrisburg · Pittsburgh · Cleveland · Akron · Wash. DC · Phil · Newark · Bohemia · Sacramento · San Francisco · San Francisco · Salt Lake City · Las Vegas · Denver · Omaha · South Bend · Dayton · Columbus · Silver Springs · Baltimore · NYC-Bdwy · Oakland · San Jose · Redwood City · Colorado Springs · Kansas City · St Louis · Florissant · Indianapolis · Cincinnati · Arlington · Camden, NJ · Honolulu · Sherman Oaks · Los Angeles · San Bernardino · Albuquerque · Louisville · Nashville · Greensboro · Charlotte · Norfolk · Richmond · Raleigh · LA-Airport Blvd · Garden a · Anaheim · Oklahoma City · Tulsa · Little Rock · Memphis · Columbia · San Diego · Phoenix · Ft. Worth · Dallas · Birmingham · Norcross · Dunwoody · Atlanta · Jacksonville · Austin · New Orleans · Orlando · San Antonio · Houston · Tampa · W. Palm Beach · Ft. Lauderdale · Ojus · Ft. Lauderdale · Miami · San Juan PR

# GARR-B

# wiscnet.net



GO BUCKY!

Legend:
- Gigabit Ethernet (purple)
- OC-12 (622Mbps) (orange)
- OC-3 (155Mbps) (blue)
- DS-3 (45Mbps) (red)
- T1 (1.5Mbps) (green)

Map locations: UW-Superior, Rice Lake, Rhinelander, UW-River Falls, UW-Stout, Marshfield, Wausau, Stiles Jct., UW-Eau Claire, Clintonville, UW-Stevens Point, Qwest and Other Provider(s), UW-Green Bay, Fox Valley TC, UW-Oshkosh, UW-La Crosse, La Crosse, Portage, Internet 2 & Qwest, Genuity, Dodgeville, UW-Madison, UW-Milwaukee, UW-Platteville, UW-Whitewater, UW-Parkside, Chicago - 1, Chicago - 2 (Winter '02)

(Summer '02), (Summer '03), (Winter '02)

Chicago
- Peering - Public and Private
- Commodity Internet Transit
- Internet2
- Merit and Other State Networks
- National Education Network
- Regional Research Peers

**ICNP 2002**

# MIT.edu

http://bgp.lcs.mit.edu/

# Network Interconnections

- **Exchange Point**
  - **Layer 2 or Layer 3**

- **Private Circuit**
  - **May be provided by a third party**

# PART II

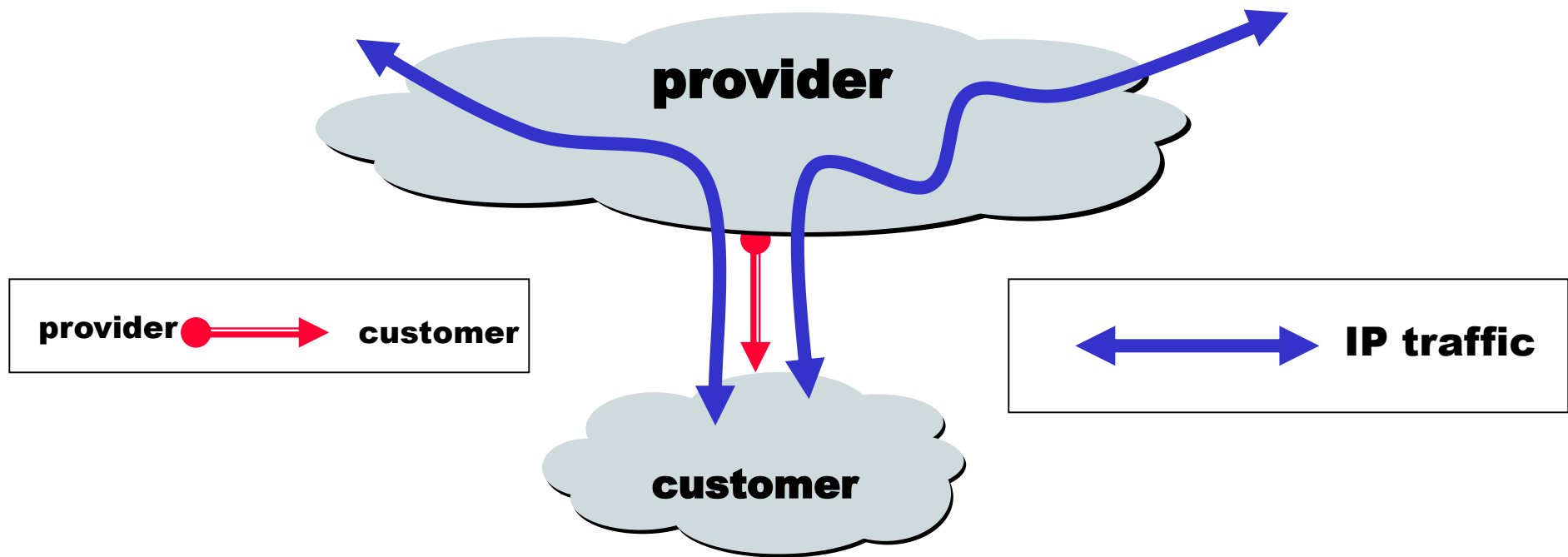## Relationships Between Networks

ICNP 2002

# Some Costs of Running an ISP

- People
- Physical connectivity and bandwidth
- Hardware
- Data center space and power
- ...

# Ballpark Figures (In US $)

- Hardware for an OC192 Pop: about 3,000,000.
  - Installation: 10,000
  - Power: 20,000/month
- OC192 link from NYC to D.C.: about 2,000,000/year
- Gigabit Ethernet IP connectivity
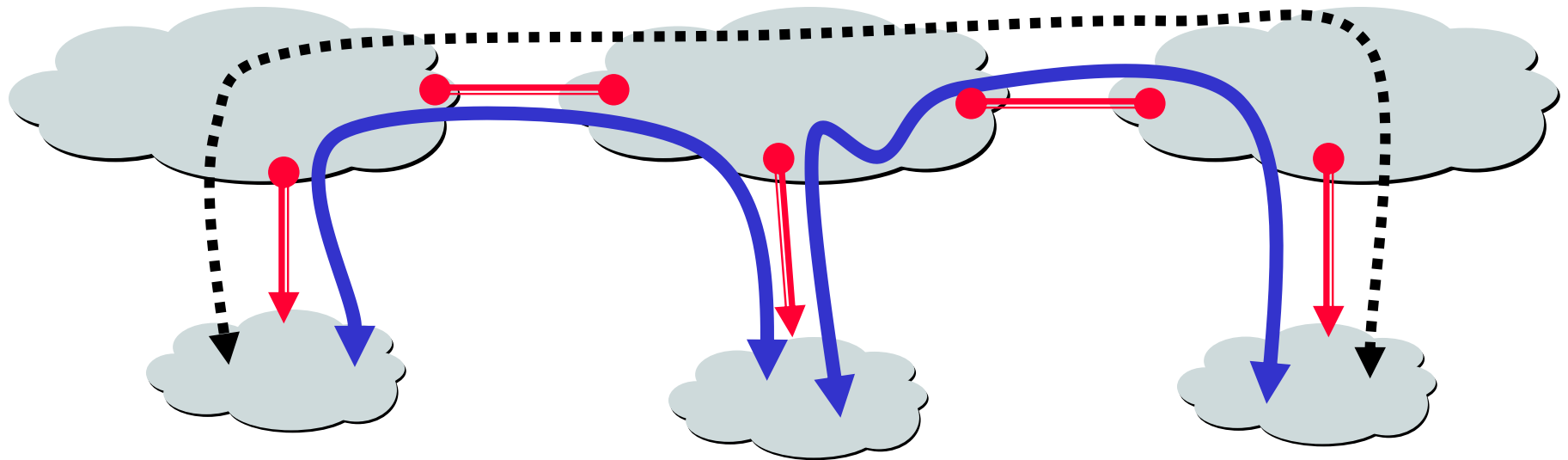  - For end user: 10,000/month
  - For ISP: 30,000/month

Prices can vary widely.  Thanks to Ben Black and Vijay Gill for hints.

# Customers and Providers



**Customer pays provider for access to the Internet**

# The "Peering" Relationship



**Legend:**
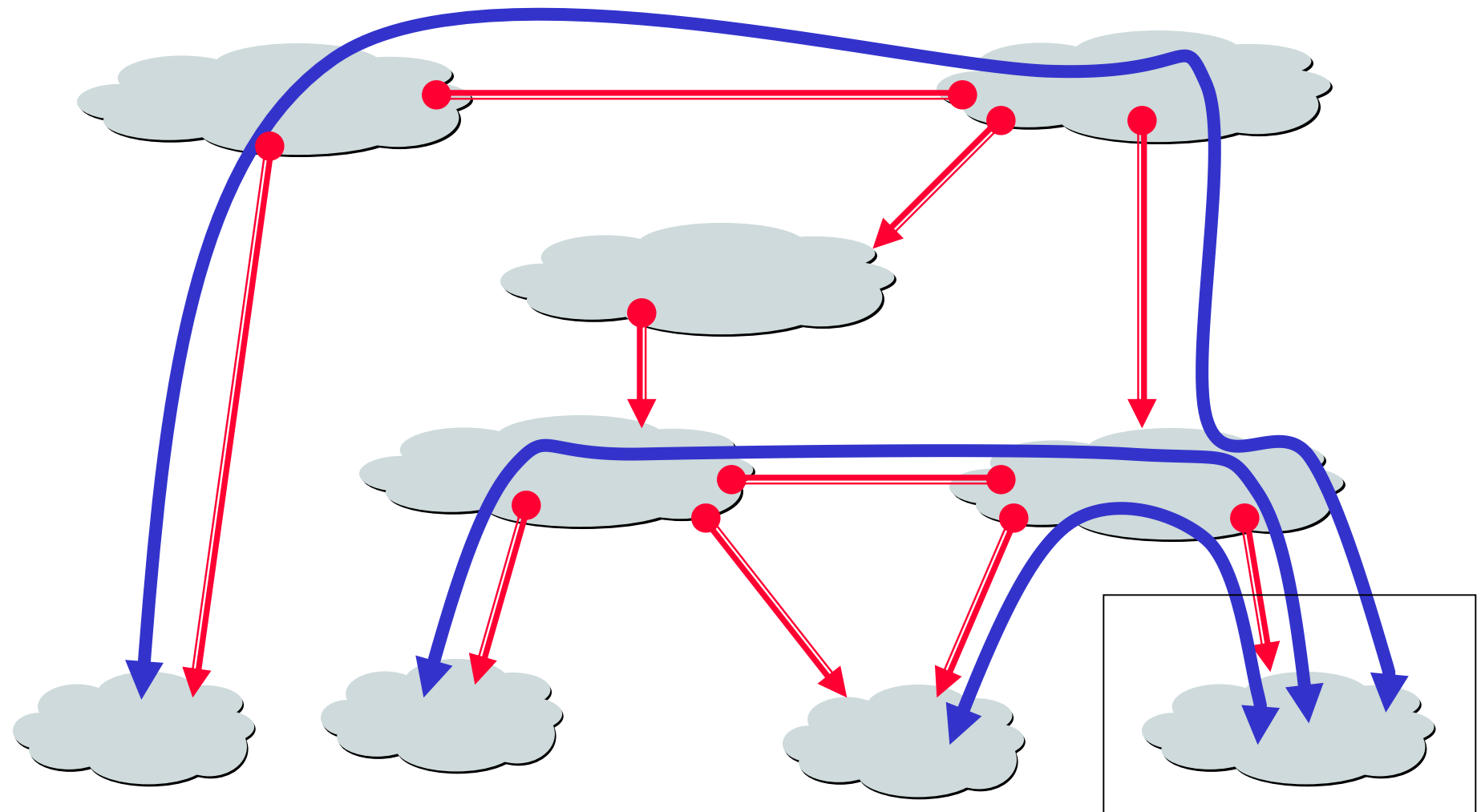
- peer ●—● peer
- provider ●—→ customer

← traffic allowed →

← traffic NOT allowed →

Peers provide transit between their respective customers
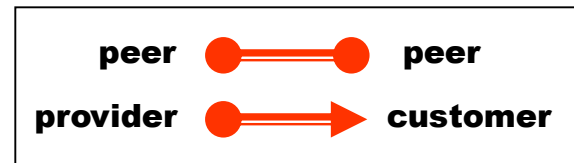
Peers do not provide transit between peers

Peers (often) do not exchange $$$

# Peering Provides Shortcuts



**Peering also allows connectivity between the customers of "Tier 1" providers.**

# Peering Wars

**Peer**

- Reduces upstream transit costs
- Can increase end-to-end performance
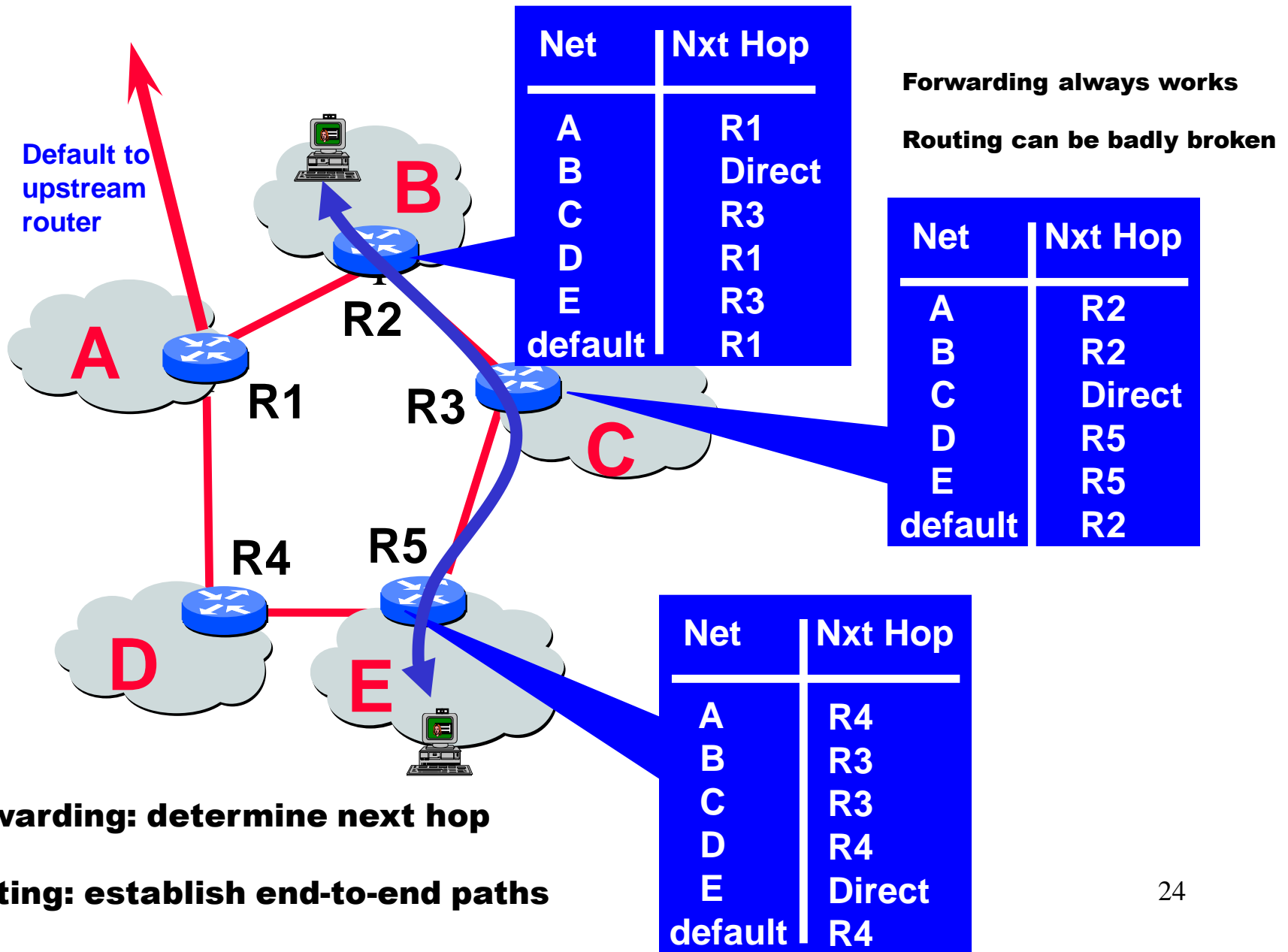- May be the only way to connect your customers to some part of the Internet ("Tier 1")

**Don't Peer**

- You would rather have customers
- Peers are usually your competition
- Peering relationships may require periodic renegotiation

Peering struggles are by far the most contentious issues in the ISP world!

Peering agreements are often confidential.

# Routing vs. Forwarding

**Default to upstream router**

| Net | Nxt Hop |
|-----|---------|
| A | R1 |
| B | Direct |
| C | R3 |
| D | R1 |
| E | R3 |
| default | R1 |

**B**

**R2**

**A**

**R1**

**R3**

**C**

**Forwarding always works**

**Routing can be badly broken**

| Net | Nxt Hop |
|-----|---------|
| A | R2 |
| B | R2 |
| C | Direct |
| D | R5 |
| E | R5 |
| default | R2 |

**R4**

**R5**

**D**

**E**

| Net | Nxt Hop |
|-----|---------|
| A | R4 |
| B | R3 |
| C | R3 |
| D | R4 |
| E | Direct |
| default | R4 |

**Forwarding: determine next hop**

**Routing: establish end-to-end paths**

24

# How Are Forwarding Tables Populated to implement Routing?

## Statically

Administrator manually configures forwarding table entries

+ More control
+ Not restricted to destination-based forwarding
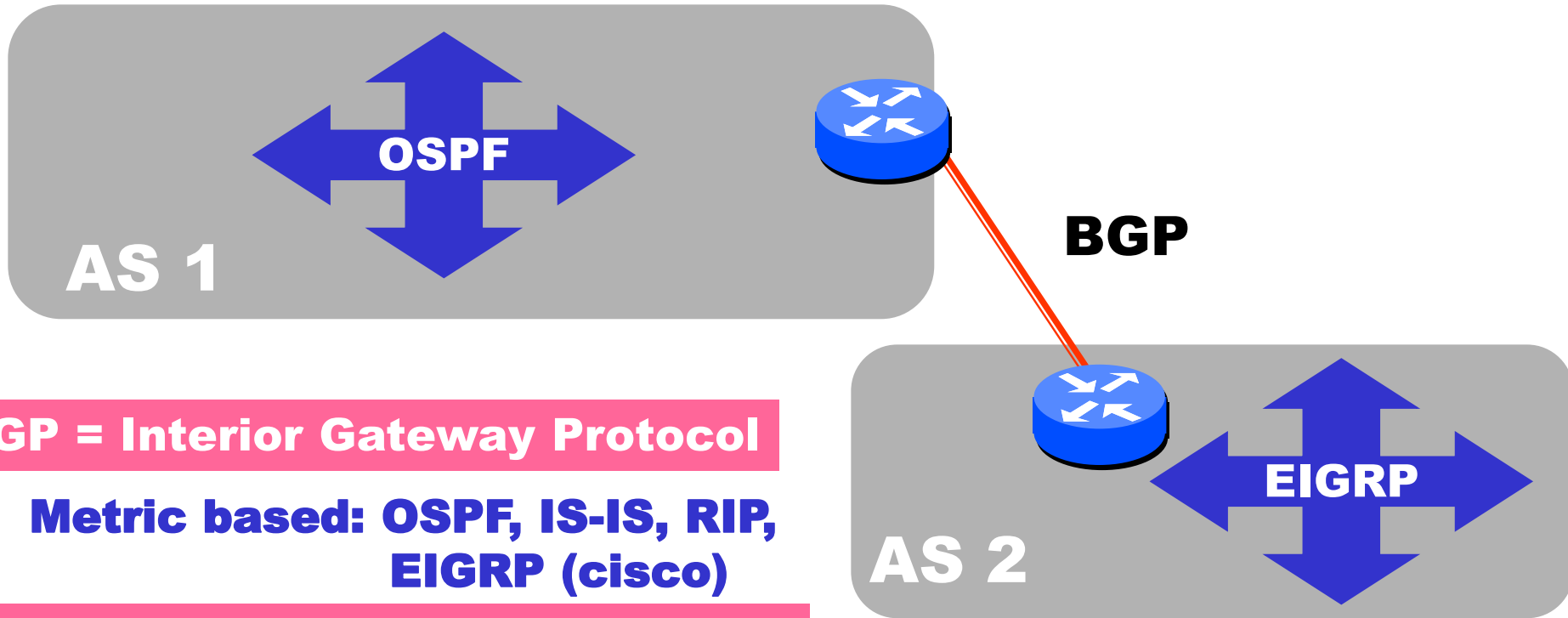- Doesn't scale
- Slow to adapt to network failures

## Dynamically

Routers exchange network reachability information using <u>ROUTING PROTOCOLS</u>. Routers use this to compute best routes

+ Can rapidly adapt to changes in network topology
+ Can be made to scale well
- Complex distributed algorithms
- Consume CPU, Bandwidth, Memory
- Debugging can be difficult
- Current protocols are destination-based

**In practice : a mix of these. Static routing mostly at the "edge"**

# Architecture of Dynamic Routing



**OSPF** — AS 1

**BGP**

**EIGRP** — AS 2

**IGP = Interior Gateway Protocol**

Metric based: OSPF, IS-IS, RIP, EIGRP (cisco)

**EGP = Exterior Gateway Protocol**

Policy based: BGP

The Routing Domain of BGP is the entire Internet

# Technology of Distributed Routing

## Link State

- Topology information is <u>flooded</u> within the routing domain
- Best end-to-end paths are computed locally at each router.
- Best end-to-end paths determine next-hops.
- Based on minimizing some notion of distance
- Works only if policy is <u>shared</u> and <u>uniform</u>
- Examples: OSPF, IS-IS

## Vectoring

- Each router knows little about network topology
- Only best next-hops are chosen by each router for each destination network.
- Best end-to-end paths result from composition of all next-hop choices
- Does not require any notion of distance
- Does not require uniform policies at all routers
- Examples: RIP, BGP

ICNP 2002

# The Gang of Four

|  | Link State | Vectoring |
|---|---|---|
| IGP | OSPF IS-IS | RIP |
| EGP |  | BGP |

ICNP 2002

# Routers Talking to Routers

**Routing info** →

← **Routing info**

- **Routing computation is distributed among routers within a routing domain**

- **Computation of best next hop based on routing information is the most CPU/memory intensive task on a router**

- **Routing messages are usually not routed, but exchanged via layer 2 between physically adjacent routers (internal BGP and multi-hop external BGP are exceptions)**

# Autonomous Routing Domains (ARDs)

A collection of physical networks glued together using IP, that have a unified administrative routing policy.

- Campus networks
- Corporate networks
- ISP Internal networks
- ...

# Autonomous Systems (ASes)

An autonomous system is an autonomous routing domain that has been assigned an Autonomous System Number (ASN).

… the administration of an AS appears to other ASes to have a single coherent interior routing plan and presents a consistent picture of what networks are reachable through it.

RFC 1930: Guidelines for creation, selection, and registration of an Autonomous System

# AS Numbers (ASNs)

ASNs are 16 bit values.
64512 through 65535 are "private"

Currently over 11,000 in use.

- **Genuity (f.k.a. BBN): 1**
- **MIT: 3**
- **Harvard: 11**
- **UC San Diego: 7377**
- **AT&T: 7018, 6341, 5074, ...**
- **UUNET: 701, 702, 284, 12199, ...**
- **Sprint: 1239, 1240, 6211, 6242, ...**
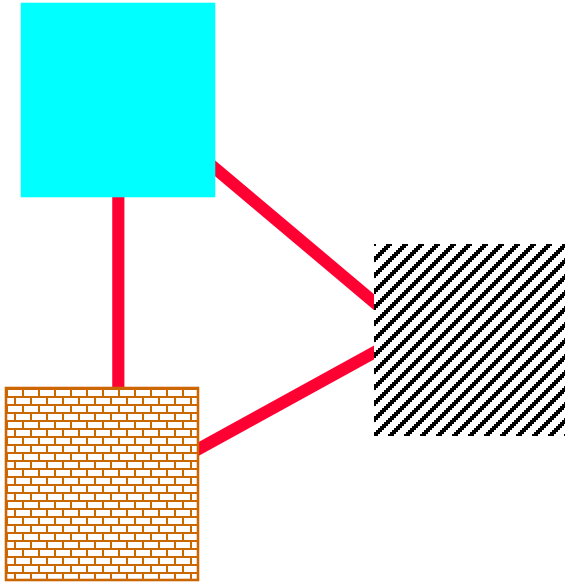- **...**

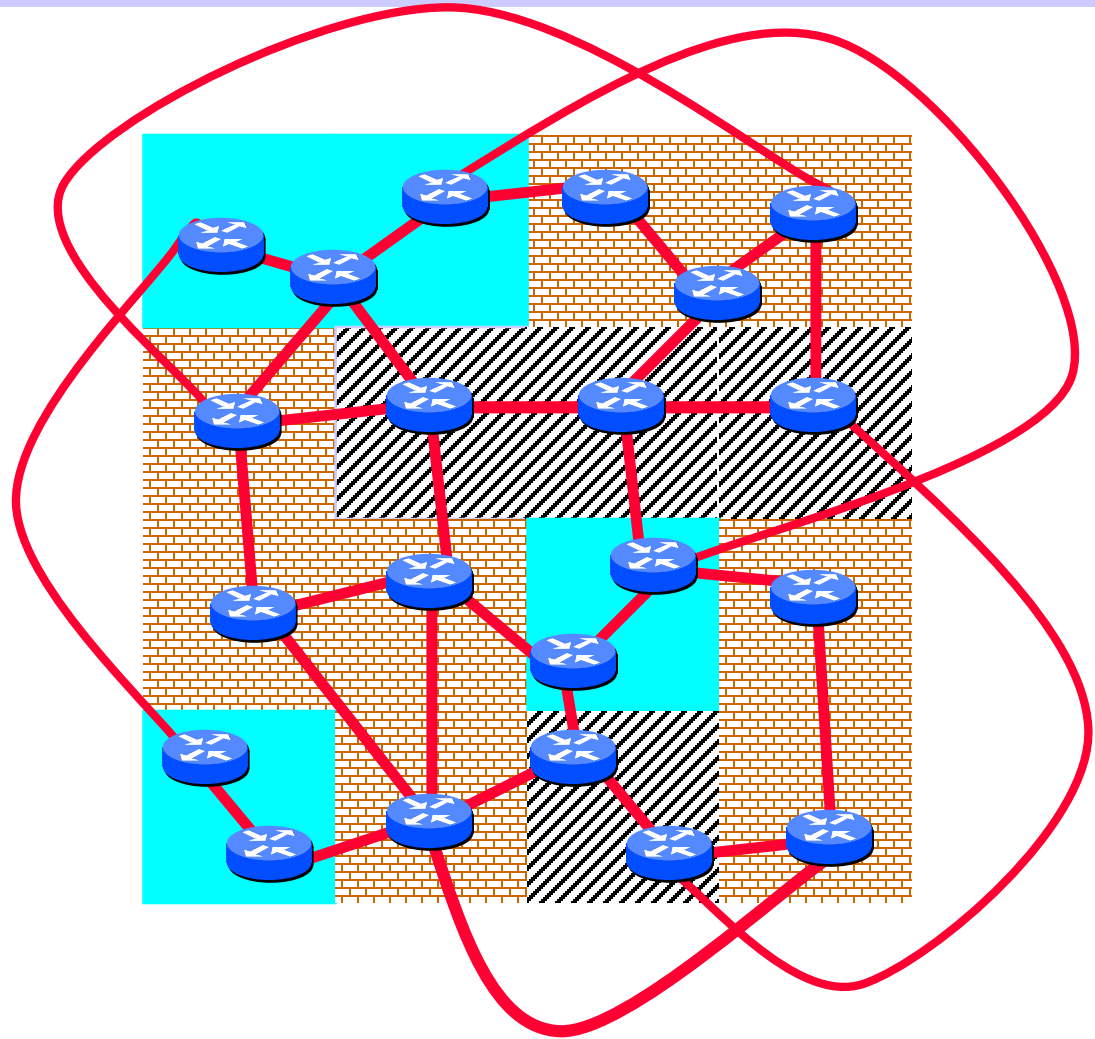**ASNs represent units of routing policy**

# AS Graphs Can Be Fun



**Part of Worldcom's Global ARD**

**AT&T North America**

701 703 1239 7018 3561 5696 6347 209 2548 3549 2914 1 702 6461 3356 174 3786 9057 6453 4766 8297 1755 3967 5511

The <u>subgraph</u> showing all ASes that have more than 100 neighbors in full graph of 11,158 nodes. July 6, 2001. **Point of view: AT&T route-server**

# AS Graph != Internet Topology
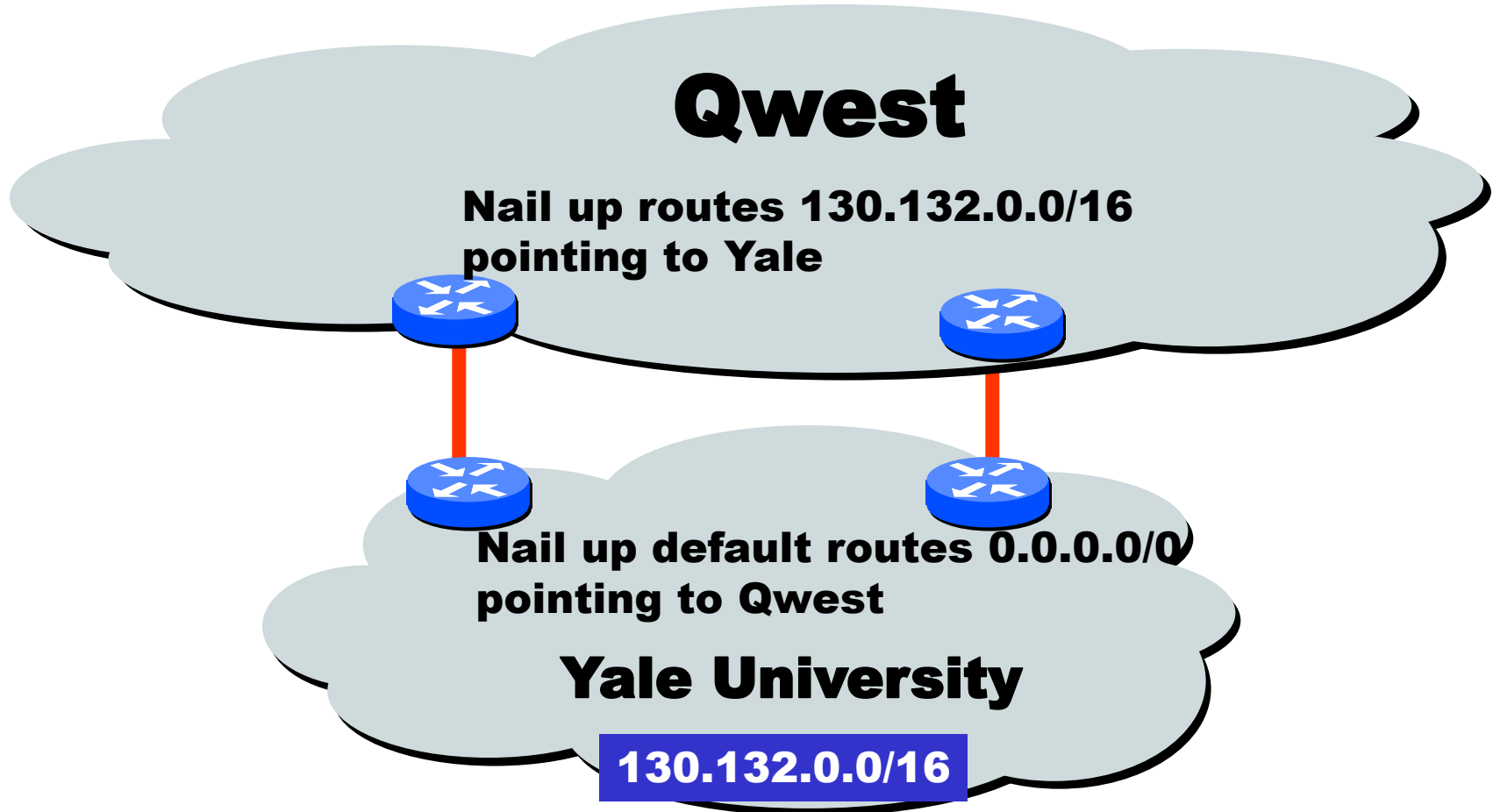
**BGP was designed to throw away information!**



**The AS graph may look like this.**
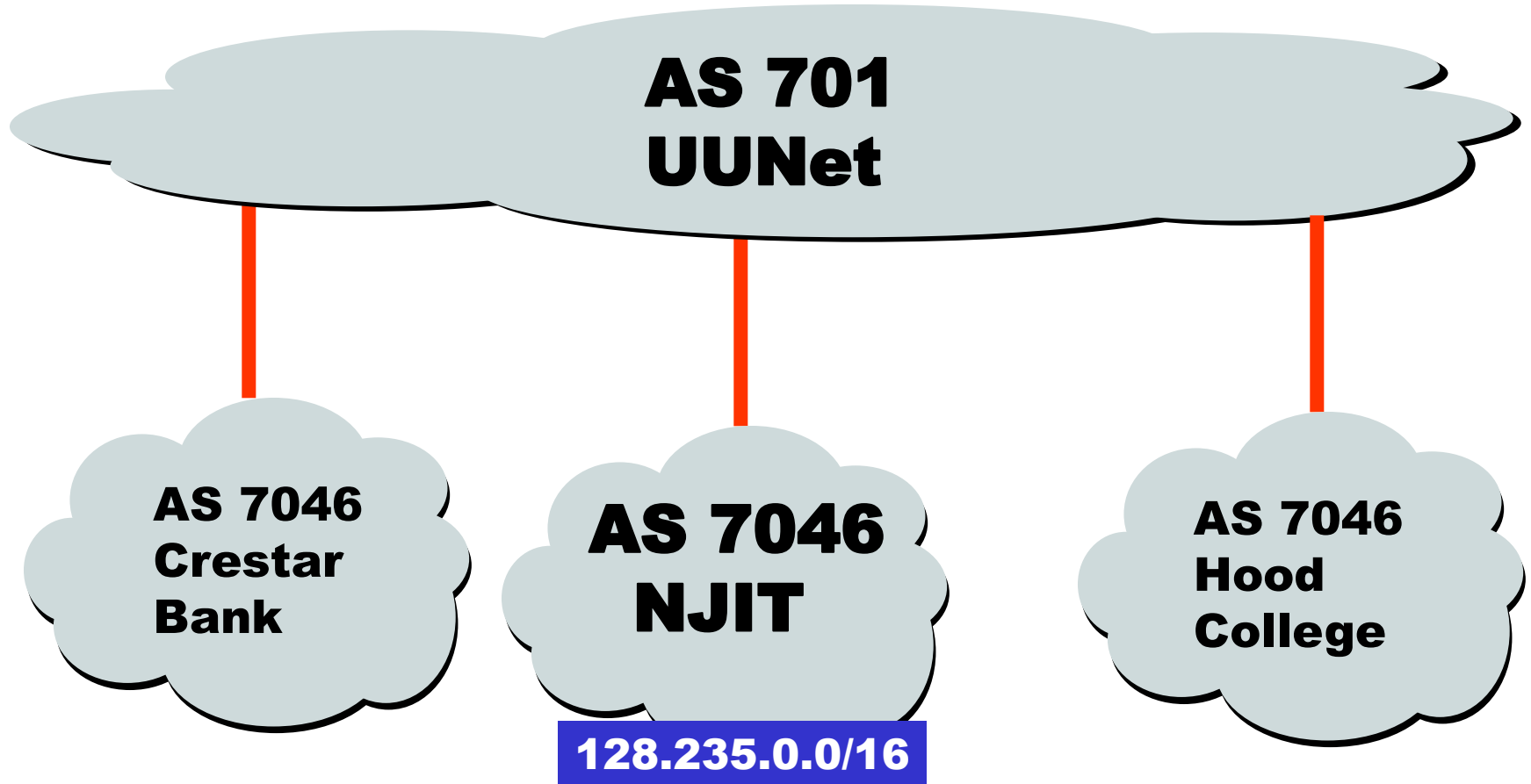
**Reality may be closer to this...**

ICNP 2002

# Autonomous Routing Domains Don't Always Need BGP or an ASN

**Qwest**

**Nail up routes 130.132.0.0/16 pointing to Yale**

**Nail up default routes 0.0.0.0/0 pointing to Qwest**

**Yale University**

**130.132.0.0/16**

**Static routing is the most common way of connecting an autonomous routing domain to the Internet.**
**This helps explain why BGP is a mystery to many ...**

ICNP 2002

# ASNs Can Be "Shared" (RFC 2270)



**AS 701 UUNet**

**AS 7046 Crestar Bank**

**AS 7046 NJIT**

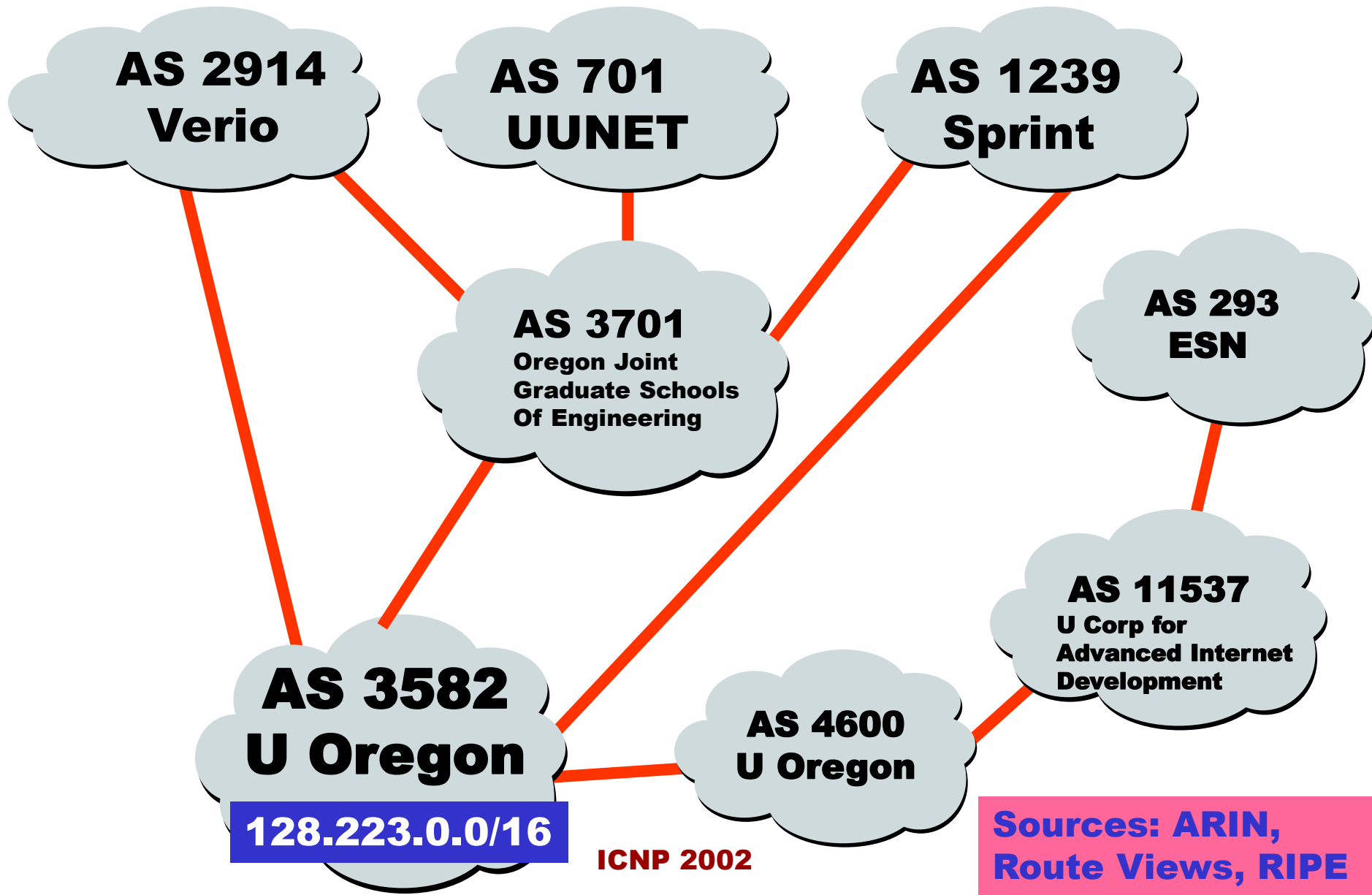**AS 7046 Hood College**

128.235.0.0/16

ASN 7046 is assigned to UUNet.  It is used by Customers single homed to UUNet, but needing BGP for some reason (load balancing, etc..) [RFC 2270]
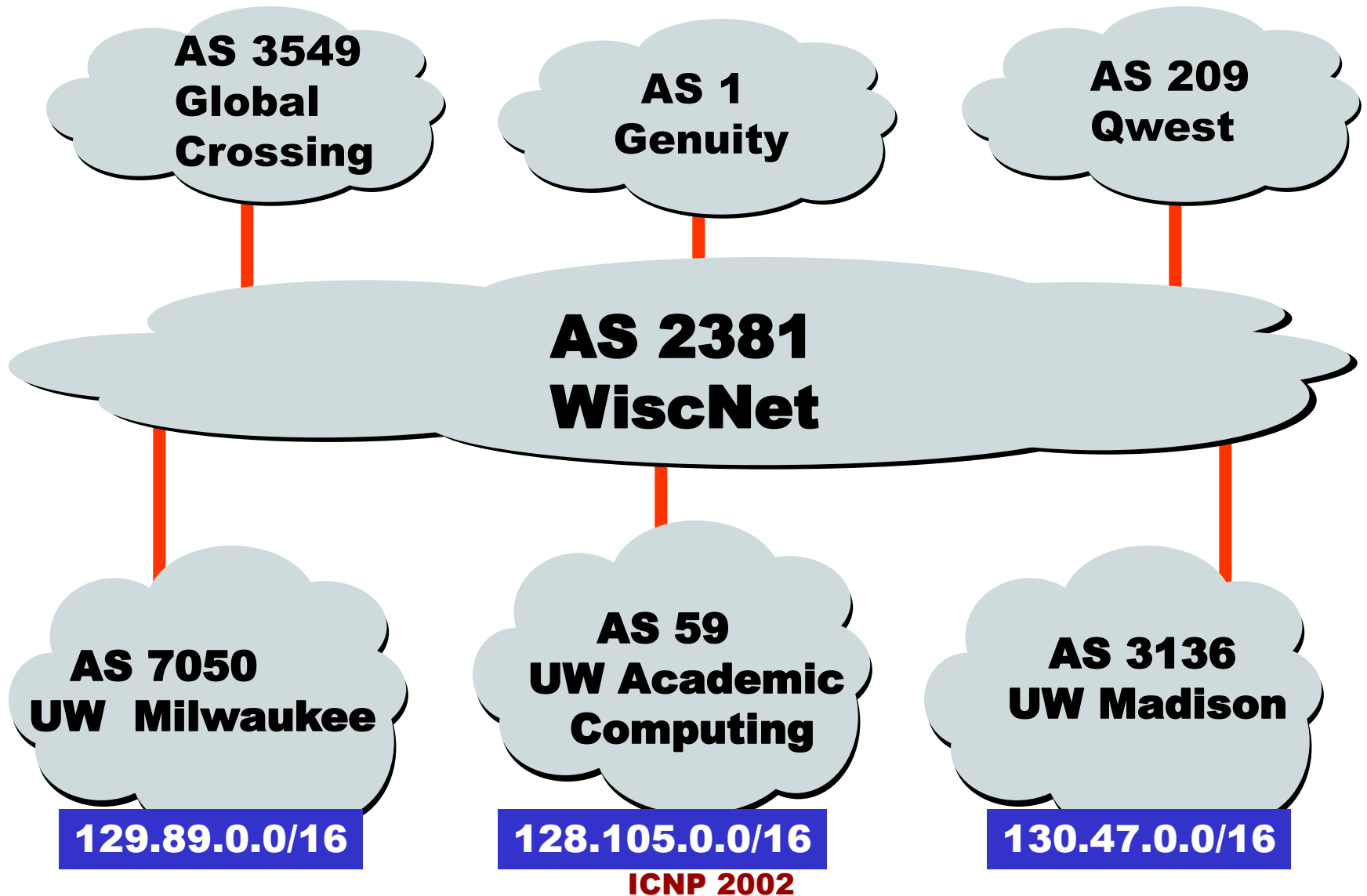
# A Bit of U Oregon's AS Neighborhood



ICNP 2002

# Partial View of cs.wisc.edu Neighborhood



AS 3549
Global
Crossing

AS 1
Genuity

AS 209
Qwest

AS 2381
WiscNet

AS 7050
UW Milwaukee

AS 59
UW Academic
Computing

AS 3136
UW Madison

129.89.0.0/16

128.105.0.0/16

130.47.0.0/16

ICNP 2002

# ARD != AS

- **Most ARDs have no ASN (statically routed at Internet edge)**

- **Some unrelated ARDs share the same ASN (RFC 2270)**

- **Some ARDs are implemented with multiple ASNs (example: Worldcom)**
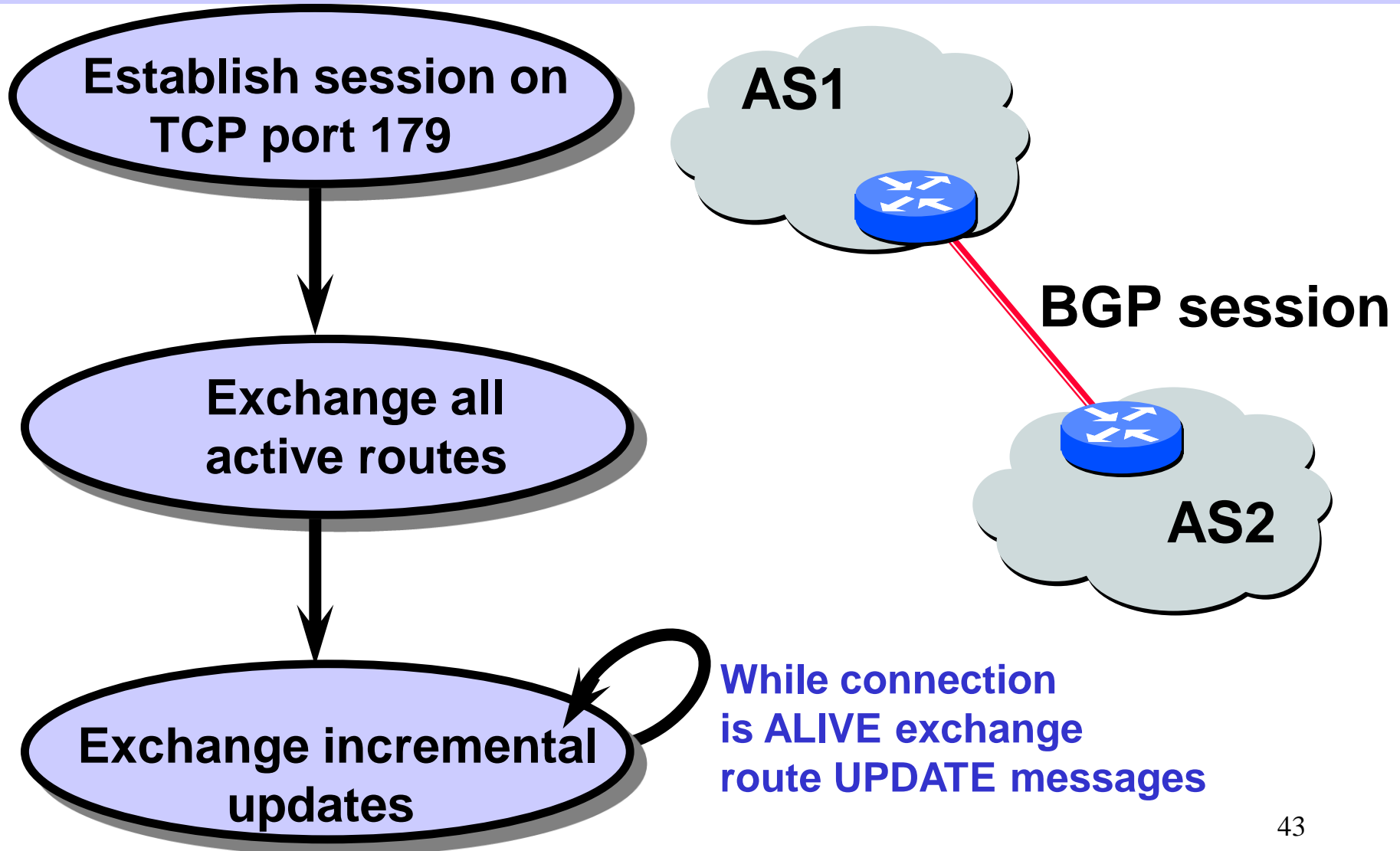
**ASes are an implementation detail of Interdomain routing**

# PART III

## Implementing Inter-Network Relationships with BGP

ICNP 2002

# BGP-4

- **BGP** = **B**order **G**ateway **P**rotocol

- Is a **Policy-Based** routing protocol

- Is the **de facto EGP** of today's global Internet

- Relatively simple protocol, but configuration is complex and the entire world can see, and be impacted by, your mistakes.
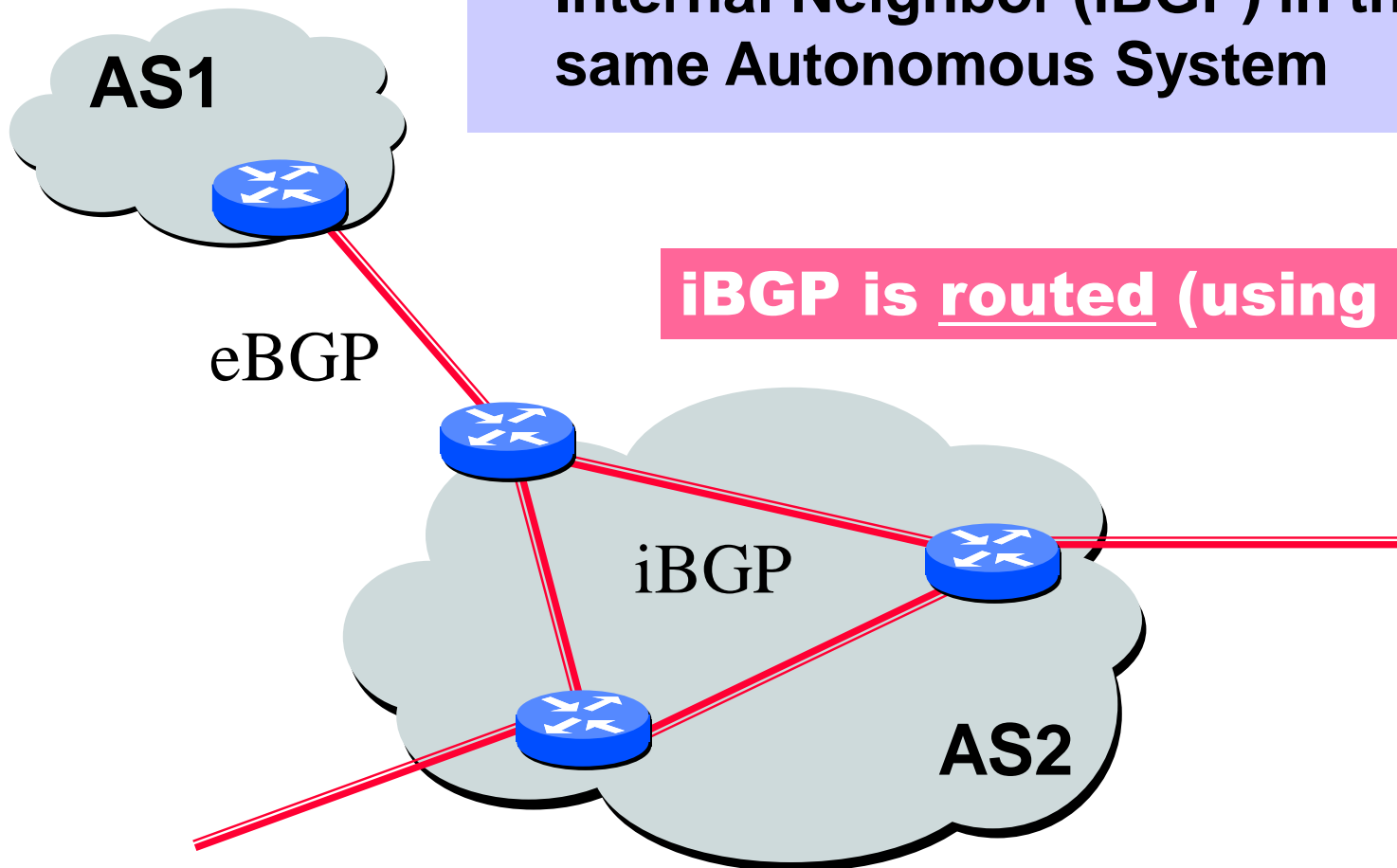
- **1989 : BGP-1 [RFC 1105]**
  - **Replacement for EGP (1984, RFC 904)**
- **1990 : BGP-2 [RFC 1163]**
- **1991 : BGP-3 [RFC 1267]**
- **1995 : BGP-4 [RFC 1771]**
  - **Support for Classless Interdomain Routing (CIDR)**

# BGP Operations (Simplified)

**Establish session on TCP port 179**

↓

**Exchange all active routes**

↓

**Exchange incremental updates** ↻

AS1

**BGP session**

AS2

**While connection is ALIVE exchange route UPDATE messages**
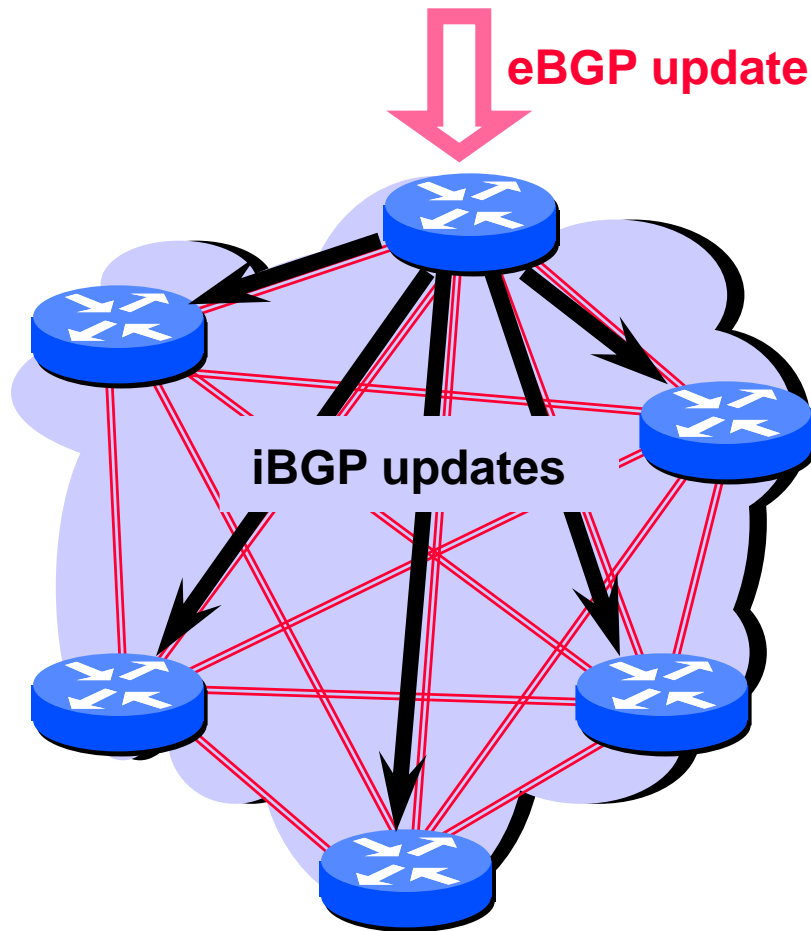
43

# Two Types of BGP Neighbor Relationships

- **External Neighbor (eBGP) in a different Autonomous Systems**
- **Internal Neighbor (iBGP) in the same Autonomous System**

**AS1**

eBGP

**iBGP is <u>routed</u> (using IGP!)**

iBGP

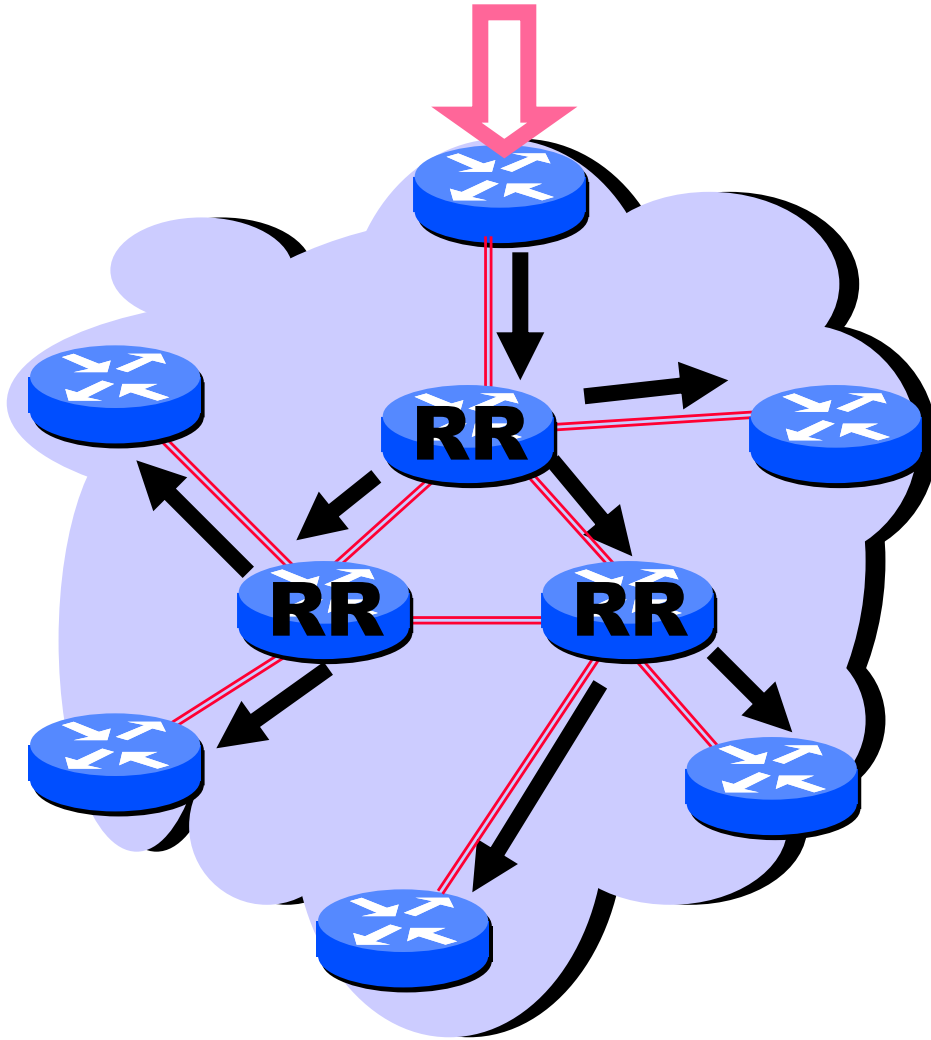**AS2**

# iBGP Mesh Does Not Scale

**eBGP update**

**iBGP updates**

- N border routers means N(N-1)/2 peering sessions

- Each router must have N-1 iBGP sessions configured

- The addition a single iBGP speaker requires configuration changes to all other iBGP speakers

- Size of iBGP routing table can be order N larger than number of best routes (remember alternate routes!)

- Each router has to listen to update noise from each neighbor

Currently four solutions:
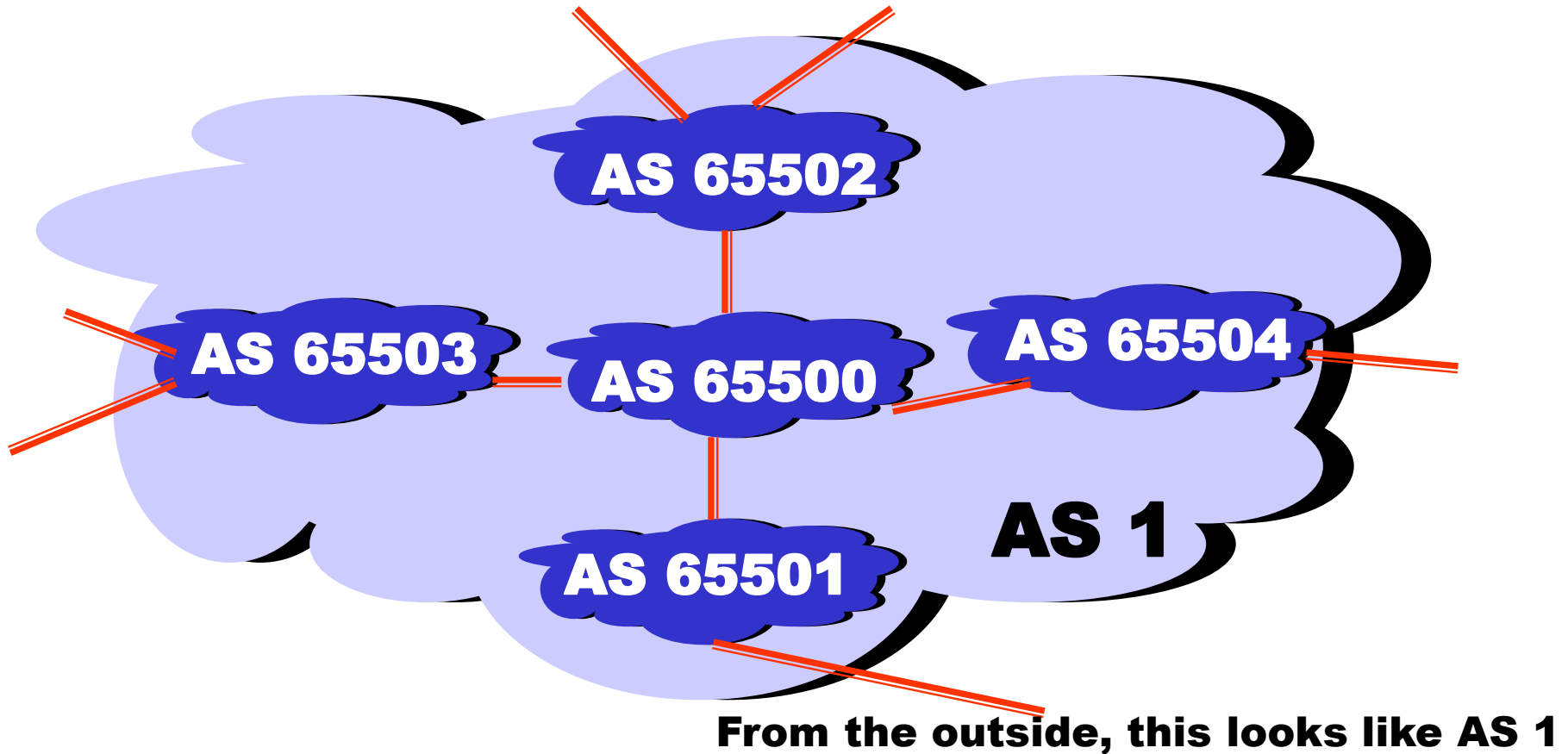(0)  Buy bigger routers!
(1)  Break AS into smaller ASes
(2)  BGP Route reflectors
(3)  BGP confederations

# Route Reflectors



- **Route reflectors can pass on iBGP updates to clients**
- **Each RR passes along ONLY best routes**
- **ORIGINATOR_ID and CLUSTER_LIST attributes are needed to avoid loops**

# BGP Confederations



AS 65502

AS 65503

AS 65500

AS 65504

AS 65501

AS 1

From the outside, this looks like AS 1

Confederation eBGP (between member ASes) preserves
LOCAL_PREF, MED, and BGP NEXTHOP.

# Four Types of BGP Messages

- **Open : Establish a peering session.**

- **Keep Alive : Handshake at regular intervals.**

- **Notification : Shuts down a peering session.**

- **Update : <u>Announcing</u> new routes or <u>withdrawing</u> previously announced routes.**

**announcement
=
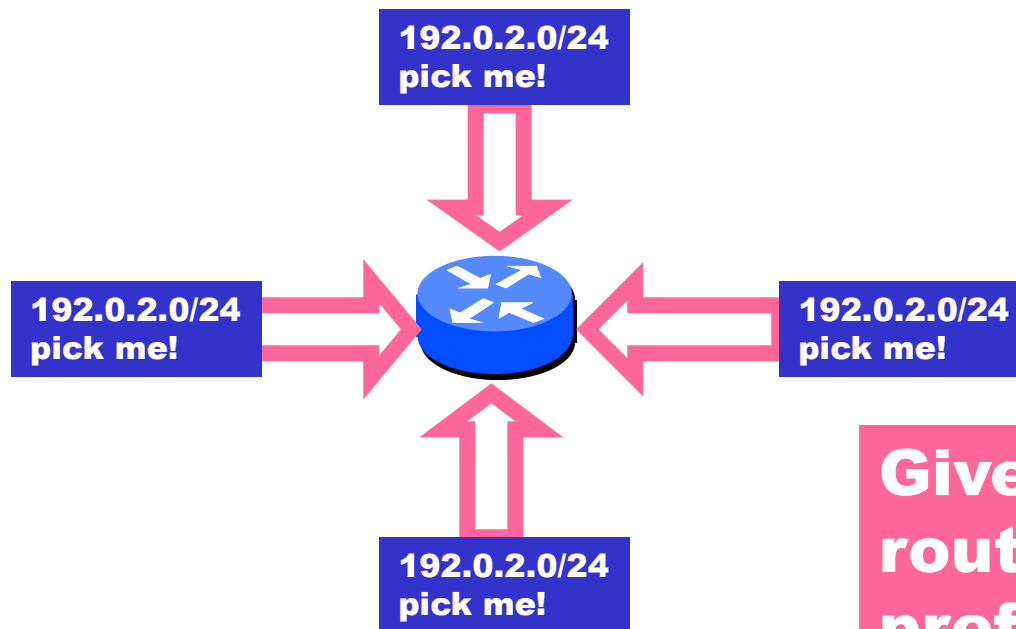prefix + <u>attributes values</u>**

# BGP Attributes

```
Value        Code                                  Reference
-----        ----------------------------------    ----------
   1         ORIGIN                                [RFC1771]
   2         AS_PATH                               [RFC1771]
   3         NEXT_HOP                              [RFC1771]
   4         MULTI_EXIT_DISC                       [RFC1771]
   5         LOCAL_PREF                            [RFC1771]
   6         ATOMIC_AGGREGATE                      [RFC1771]
   7         AGGREGATOR                            [RFC1771]
   8         COMMUNITY                             [RFC1997]
   9         ORIGINATOR_ID                         [RFC2796]
  10         CLUSTER_LIST                          [RFC2796]
  11         DPA                                      [Chen]
  12         ADVERTISER                            [RFC1863]
  13         RCID_PATH / CLUSTER_ID                [RFC1863]
  14         MP_REACH_NLRI                         [RFC2283]
  15         MP_UNREACH_NLRI                       [RFC2283]
  16         EXTENDED COMMUNITIES                    [Rosen]
...
 255         reserved for development
```

Most important attributes

From IANA: http://www.iana.org/assignments/bgp-parameters

Not all attributes need to be present in every announcement

ICNP 2002

# Attributes are Used to Select Best Routes

192.0.2.0/24
pick me!

192.0.2.0/24
pick me!

192.0.2.0/24
pick me!

192.0.2.0/24
pick me!

Given multiple routes to the same prefix, a BGP speaker must pick at most <u>one</u> best route

(Note: it could reject them all!)

# Route Selection Summary

**Highest Local Preference**     Enforce relationships
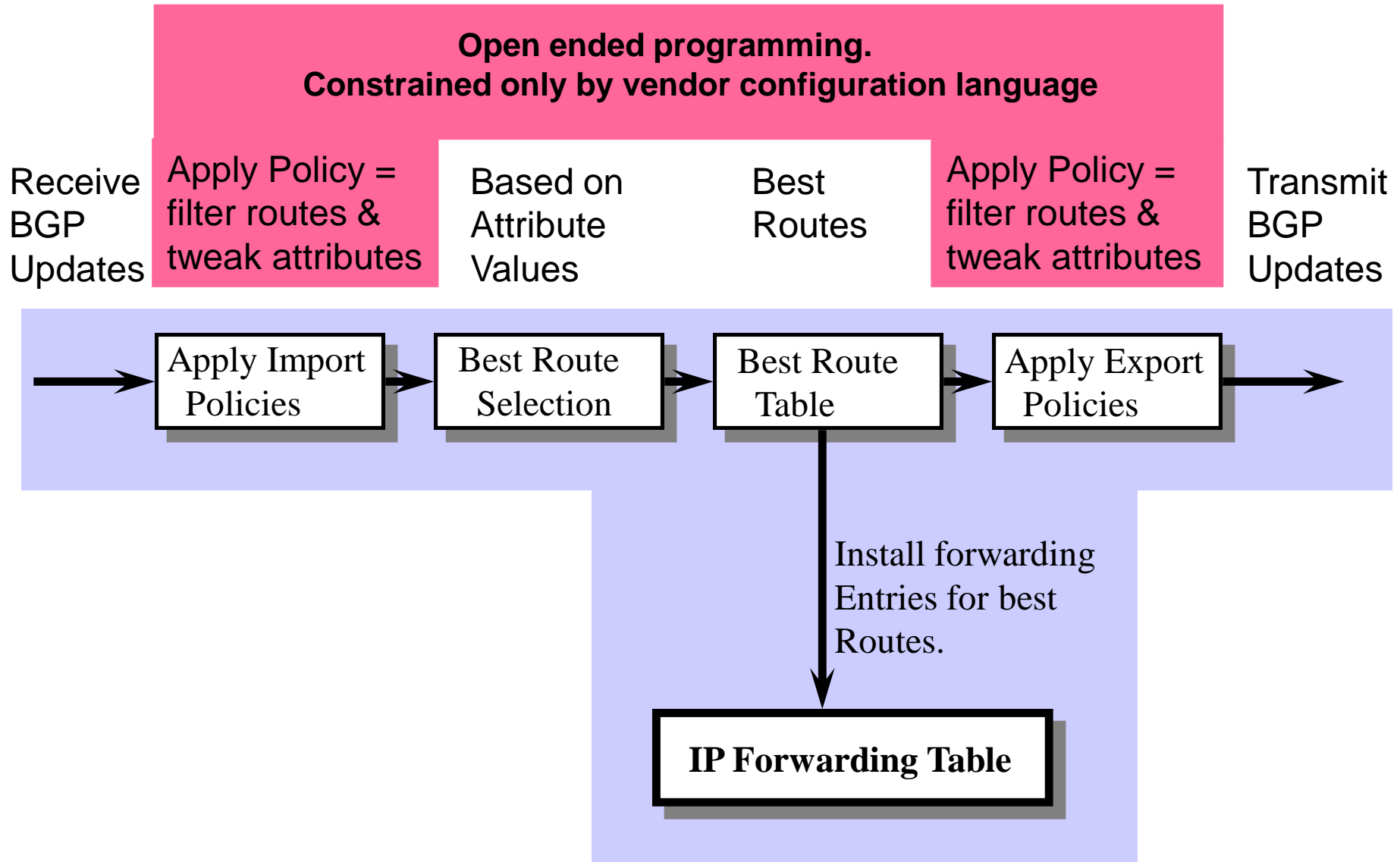
**Shortest ASPATH**

**Lowest MED**

**i-BGP < e-BGP**                 traffic engineering

**Lowest IGP cost**
**to BGP egress**

**Lowest router ID**             Throw up hands and
                                 break ties

# BGP Route Processing

Open ended programming.
Constrained only by vendor configuration language

| Receive BGP Updates | Apply Policy = filter routes & tweak attributes | Based on Attribute Values | Best Routes | Apply Policy = filter routes & tweak attributes | Transmit BGP Updates |

→ **Apply Import Policies** → **Best Route Selection** → **Best Route Table** → **Apply Export Policies** →

Install forwarding Entries for best Routes.

**IP Forwarding Table**

# BGP Next Hop Attribute



**12.125.133.90**

**AS 7018**
**AT&T**

**12.127.0.121**

**AS 6431**
**AT&T Research**

**AS 12654**
**RIPE NCC**
**RIS project**

**135.207.0.0/16**
**Next Hop = 12.125.133.90**

**135.207.0.0/16**
**Next Hop = 12.127.0.121**

**Every time a route announcement crosses an AS boundary, the Next Hop attribute is changed to the IP address of the border router that announced the route.**

53

# Join EGP with IGP For Connectivity

135.207.0.0/16
Next Hop = 192.0.2.1

135.207.0.0/16

10.10.10.10

AS 1

192.0.2.1

192.0.2.0/30

AS 2

**Forwarding Table**

| destination | next hop |
|---|---|
| 192.0.2.0/30 | 10.10.10.10 |

**+**

**EGP**

| destination | next hop |
|---|---|
| 135.207.0.0/16 | 192.0.2.1 |

**Forwarding Table**

| destination | next hop |
|---|---|
| 135.207.0.0/16 | 10.10.10.10 |
| 192.0.2.0/30 | 10.10.10.10 |

ICNP 2002

# Implementing Customer/Provider and Peer/Peer relationships

## Two parts:

- **Enforce  transit relationships**
  - Outbound route filtering
- **Enforce order of route preference**
  - provider < peer < customer

# Import Routes

# Export Routes



ICNP 2002

# How Can Routes be Colored? BGP Communities!

**A community value is 32 bits**

By convention, first 16 bits is ASN indicating who is giving it an interpretation

community number

Used for signally within and between ASes

Very powerful BECAUSE it has no (predefined) meaning

**Community Attribute = a list of community values. (So one route can belong to multiple communities)**

RFC 1997 (August 1996)

**Two reserved communities**
no_export = 0xFFFFFF01: don't export out of AS

no_advertise 0xFFFFFF02: don't pass to BGP neighbors

# Communities Example

**Import**

- **1:100** ♥
  - Customer routes
- **1:200** ✚
  - Peer routes
- **1:300** ◆
  - Provider Routes

**Export**

- **To Customers**
  - 1:100, 1:200, 1:300
- **To Peers**
  - 1:100
- **To Providers**
  - 1:100

**AS 1**

# So Many Choices

peer ⬤━━⬤ peer
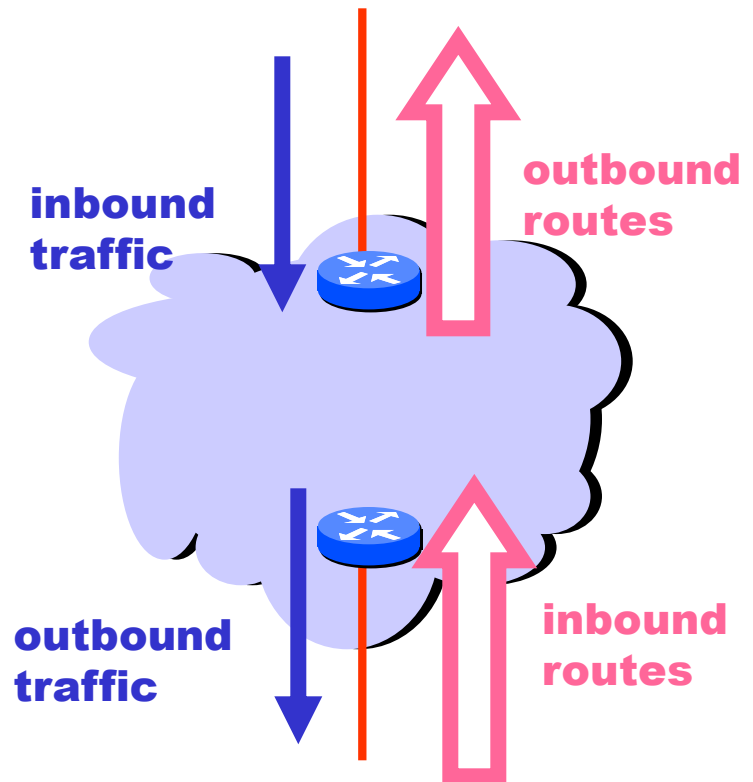provider ⬤━➡ customer

AS 4

Frank's Internet Barn

AS 3

AS 2

AS 1

13.13.0.0/16

Which route should Frank pick to 13.13.0.0./16?

# LOCAL PREFERENCE

**Local preference used ONLY in iBGP**

AS 4

local pref = 80

local pref = 90

AS 3

local pref = 100

AS 2

AS 1

**13.13.0.0/16**

**Higher Local preference values are more preferred**

61

# PART IV

## Traffic Engineering with BGP

# Tweak Tweak Tweak

- **For <u>inbound</u> traffic**
  - Filter outbound routes
  - Tweak attributes on <u>outbound</u> routes in the hope of influencing your neighbor's best route selection

- **For <u>outbound</u> traffic**
  - Filter <u>inbound</u> routes
  - Tweak attributes on <u>inbound</u> routes to influence best route selection

inbound traffic

outbound routes

outbound traffic

inbound routes

**In general, an AS has more control over outbound traffic**

# ASPATH Attribute

**135.207.0.0/16**
**AS Path = 1755 1239 7018 6341**

**AS 1129**
Global Access

**AS 1755**
Ebone

**135.207.0.0/16**
**AS Path = 1239 7018 6341**

**135.207.0.0/16**
**AS Path = 1129 1755 1239 7018 6341**

**AS 1239**
Sprint

**135.207.0.0/16**
**AS Path = 7018 6341**

**AS 12654**
RIPE NCC
RIS project

**AS7018**
AT&T

**135.207.0.0/16**
**AS Path = 3549 7018 6341**

**135.207.0.0/16**
**AS Path = 6341**

**AS 6341**
AT&T Research

**135.207.0.0/16**
**AS Path = 7018 6341**

**AS 3549**
Global Crossing

**135.207.0.0/16**

Prefix Originated

64

# Shorter Doesn't Always Mean Shorter

**Mr. BGP says that path 4 1 is better than path 3 2 1**

**Duh!**

**In fairness: could you do this "right" and still scale?**

**Exporting internal state would dramatically increase global instability and amount of routing state**

AS 4

AS 3

AS 2

AS 1

# Interdomain Loop Prevention

**BGP at AS YYY will never accept a route with ASPATH containing YYY.**

AS 7018

**Don't Accept!**

12.22.0.0/16
ASPATH = 1 333 7018 877

AS 1

# Traffic Often Follows ASPATH

135.207.0.0/16
ASPATH = 3 2 1

**AS 1**   **AS 2**   **AS 3**   **AS 4**

135.207.0.0/16

IP Packet
Dest =
135.207.44.66

# ... But It Might Not

AS 2 filters all subnets with masks longer than /24

135.207.0.0/16
ASPATH = 1

135.207.44.0/25
ASPATH = 5

135.207.0.0/16
ASPATH = 3 2 1

**AS 1**

135.207.0.0/16

**AS 2**

**AS 3**

**AS 4**

**AS 5**

135.207.44.0/25

IP Packet
Dest =
135.207.44.66

From AS 4, it may look like this packet will take path 3 2 1, but it actually takes path 3 2 5
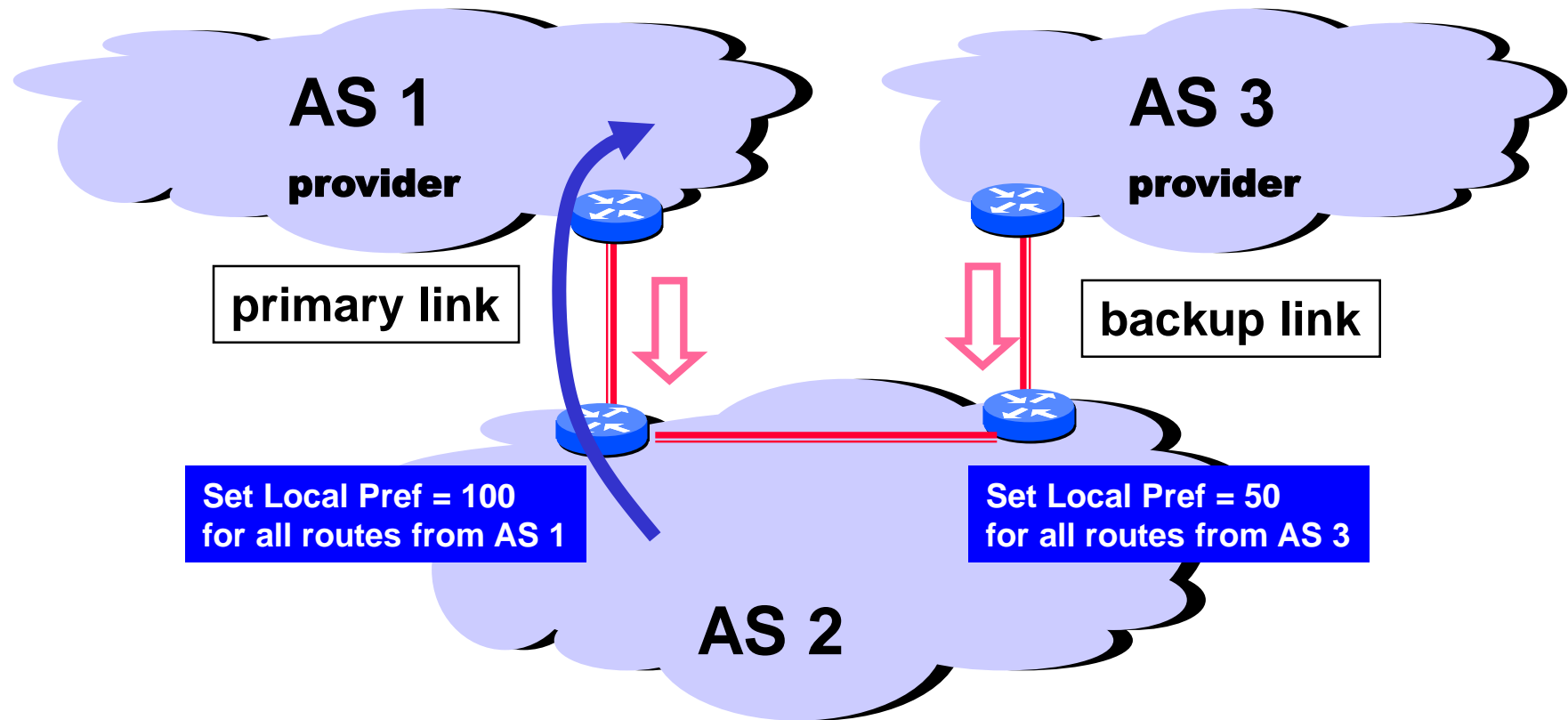
# AS Graphs Depend on Point of View



ICNP 2002

# Implementing Backup Links with Local Preference (Outbound Traffic)

AS 1

**primary link**

**backup link**

**Set Local Pref = 100 for all routes from AS 1**

**Set Local Pref = 50 for all routes from AS 1**

AS 65000

**Forces <u>outbound</u> traffic to take primary link, unless link is down.**

**We'll talk about <u>inbound</u> traffic soon ...**

# Multihomed Backups (Outbound Traffic)



AS 1
provider

AS 3
provider

primary link

backup link

Set Local Pref = 100
for all routes from AS 1

Set Local Pref = 50
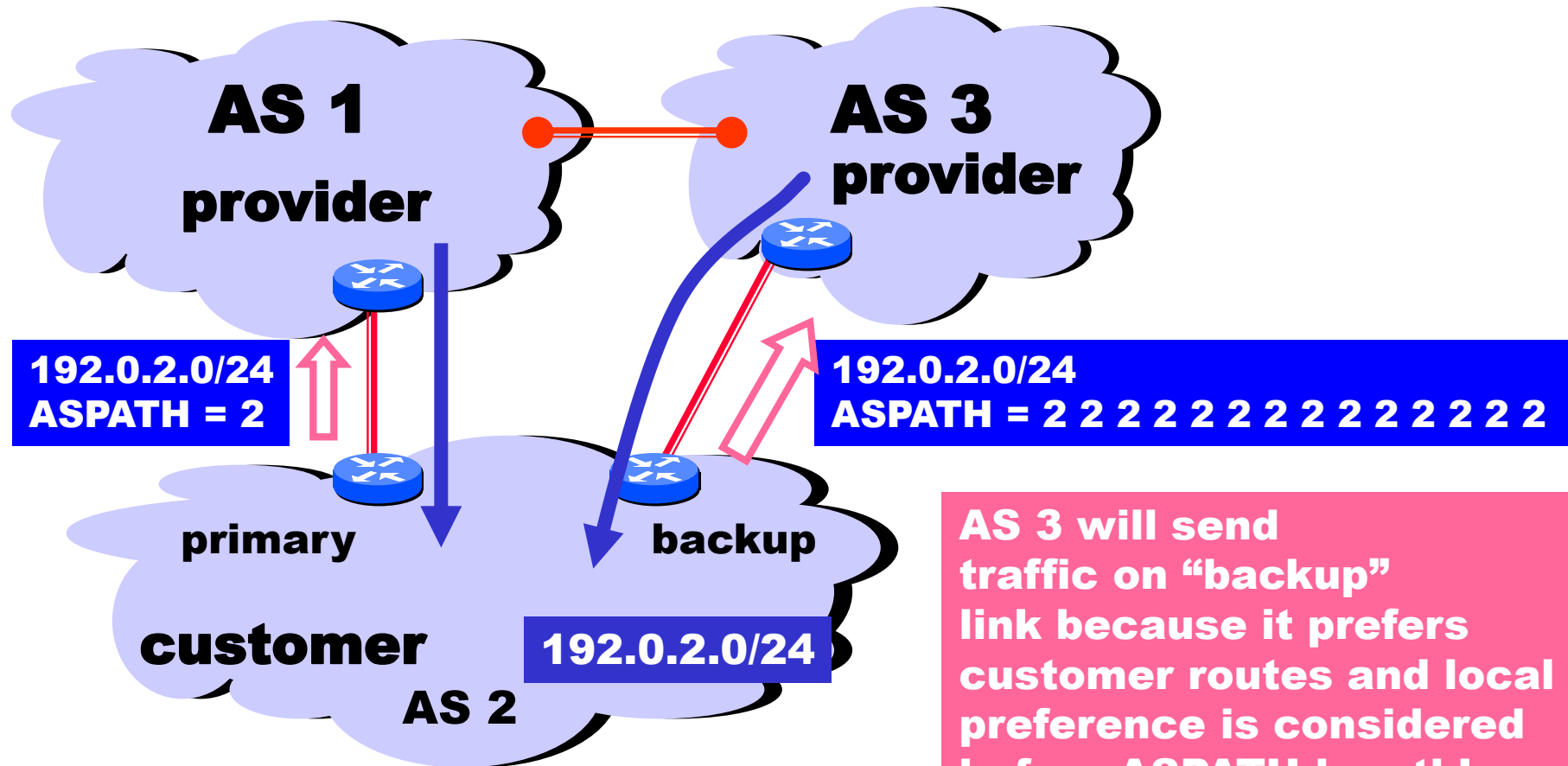for all routes from AS 3

AS 2

**Forces <u>outbound</u> traffic to take primary link, unless link is down.**

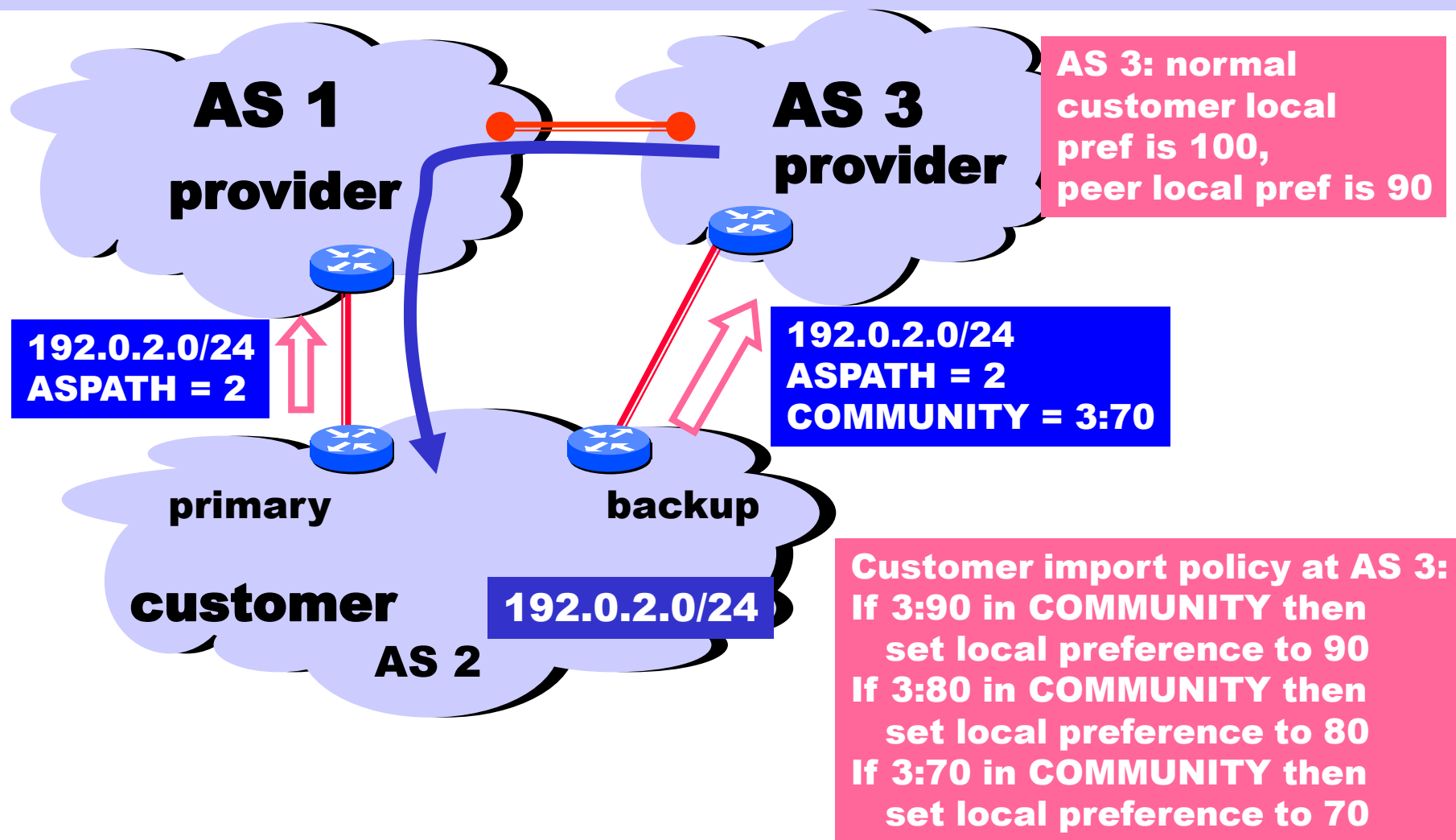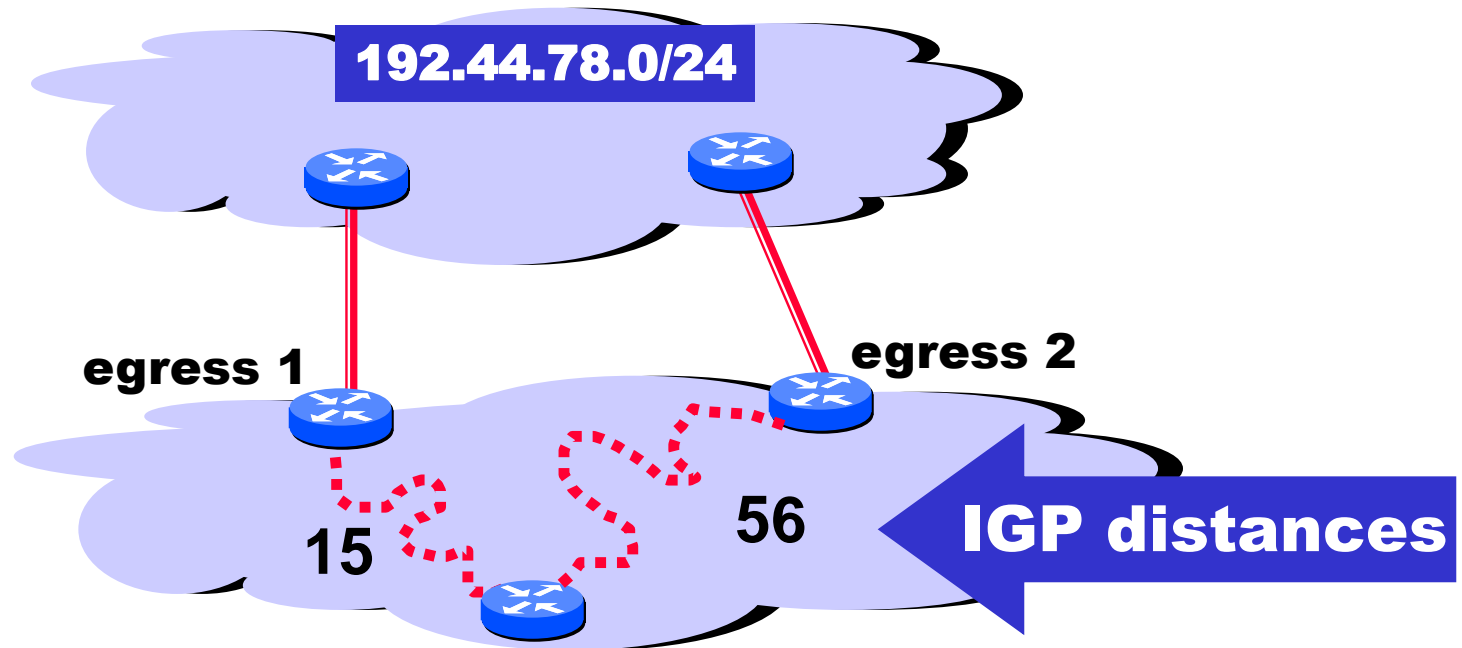# Shedding Inbound Traffic with ASPATH Padding.  Yes, this is a Glorious Hack ...



AS 1       provider

192.0.2.0/24
ASPATH = 2

192.0.2.0/24
ASPATH = 2  2  2

primary            backup

customer       192.0.2.0/24

AS 2

Padding will (usually) force inbound traffic from AS 1 to take primary link

# ... But Padding Does Not Always Work

AS 1
provider

AS 3
provider

192.0.2.0/24
ASPATH = 2

192.0.2.0/24
ASPATH = 2 2 2 2 2 2 2 2 2 2 2 2 2

primary

backup

customer
AS 2

192.0.2.0/24

AS 3 will send traffic on "backup" link because it prefers customer routes and local preference is considered before ASPATH length!

Padding in this way is often used as a form of load balancing
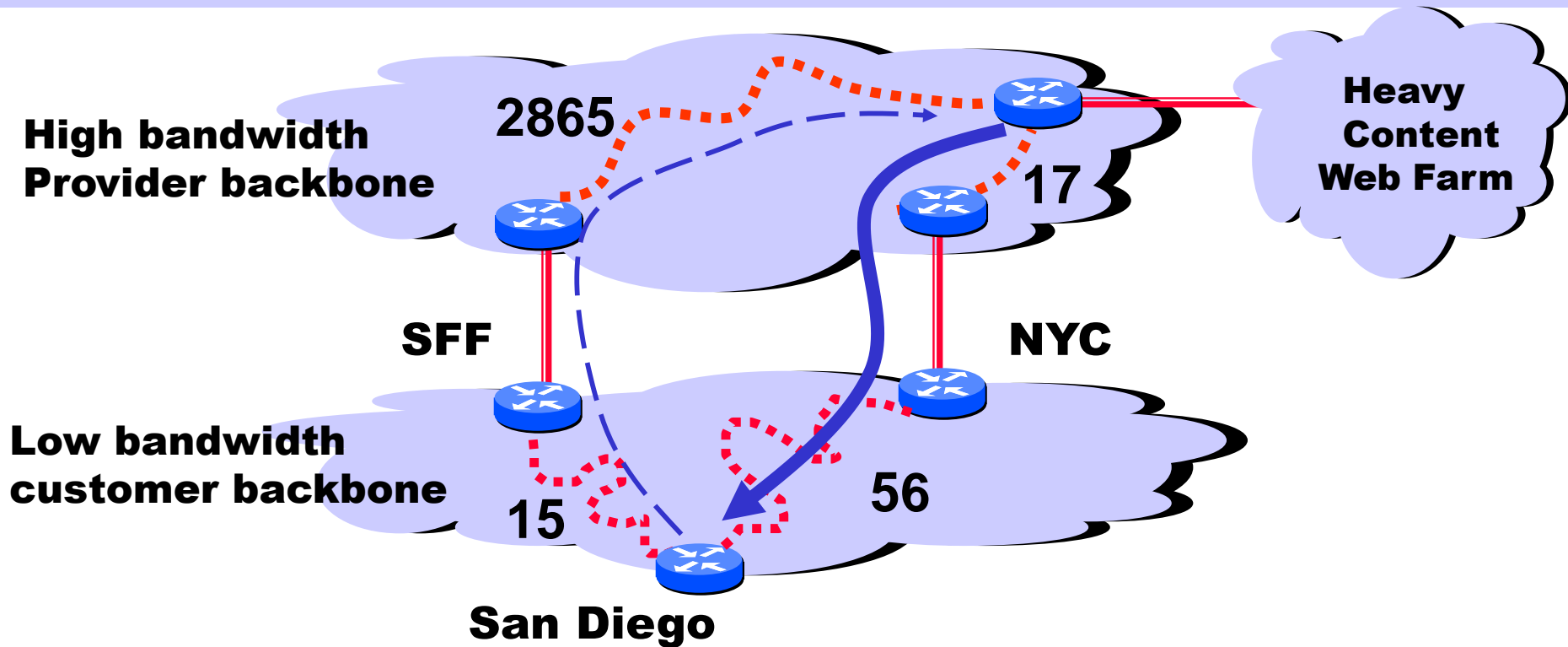
# COMMUNITY Attribute to the Rescue!

AS 1 provider

AS 3 provider

AS 3: normal customer local pref is 100, peer local pref is 90

192.0.2.0/24
ASPATH = 2

192.0.2.0/24
ASPATH = 2
COMMUNITY = 3:70

primary

backup

customer

192.0.2.0/24

AS 2

Customer import policy at AS 3:
If 3:90 in COMMUNITY then
    set local preference to 90
If 3:80 in COMMUNITY then
    set local preference to 80
If 3:70 in COMMUNITY then
    set local preference to 70

# Hot Potato Routing: Go for the Closest Egress Point



192.44.78.0/24

egress 1

egress 2

15

56

IGP distances

This Router has two BGP routes to 192.44.78.0/24.

Hot potato: get traffic off of your network as Soon as possible. Go for egress 1!

# Getting Burned by the Hot Potato



**High bandwidth Provider backbone**

2865

**Heavy Content Web Farm**

17

**SFF**

**NYC**

**Low bandwidth customer backbone**

15

56

**San Diego**

**Many customers want their provider to carry the bits!**

- - - - - → tiny http request
━━━━━━➤ huge http reply

# Cold Potato Routing with MEDs
## (Multi-Exit Discriminator Attribute)



**Prefer lower MED values**

2865

17

**Heavy Content Web Farm**

192.44.78.0/24
MED = 15

192.44.78.0/24
MED = 56

15

56

192.44.78.0/24

**This means that MEDs must be considered BEFORE IGP distance!**

**Note1 : some providers will not listen to MEDs**

**Note2 : MEDs need not be tied to IGP distance**

# PART V

## A Wee Bit O' Theory: What Problem is BGP Attempting to Solve?

# Policies Can Interact Strangely ("Route Pinning" Example)



**1**

**2** Install backup link using community

**3** Disaster strikes primary link and the backup takes over

**4** Primary link is restored but some traffic remains *pinned* to backup

ICNI 2002

# News at 11:00h

- **BGP <u>is not guaranteed</u> to converge on a stable routing. Policy interactions could lead to "livelock" protocol oscillations.**
  **See "Persistent Route Oscillations in Inter-domain Routing" by K. Varadhan, R. Govindan, and D. Estrin. ISI report, 1996**

- **Corollary: BGP <u>is not guaranteed</u> to recover from network failures.**

# What Problem is BGP Solving?

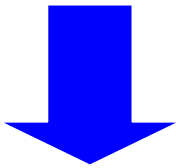| Underlying problem | Distributed means of computing a solution. |
|---|---|
| **Shortest Paths** | **RIP, OSPF, IS-IS** |
| **X?** | **BGP** |

# Separate dynamic and static semantics

static
semantics

dynamic
semantics

**BGP Policies** **BGP**

Booo Hooo,
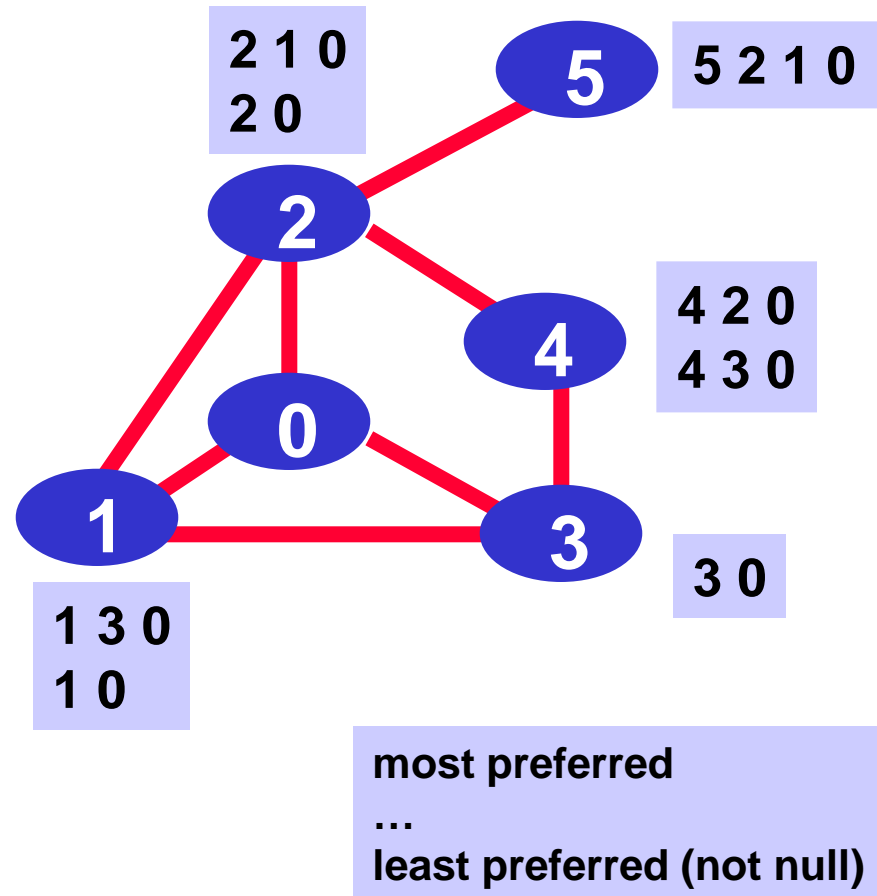Many, many
complications...

**Stable Paths
Problem (SPP)** **SPVP**

SPVP = Simple Path
Vector Protocol = a
distributed
algorithm for
solving SPP

**See [Griffin, Shepherd, Wilfong]**

# An instance of the *Stable Paths Problem* (SPP)

**2 1 0**
**2 0**

**5**

**5 2 1 0**

**2**

- A graph of nodes and edges,
- Node 0, called *the origin*,
- For each non-zero node, a set or permitted paths to the origin. This set always contains the "null path".
- A ranking of permitted paths at each node. Null path is always least preferred. (Not shown in diagram)
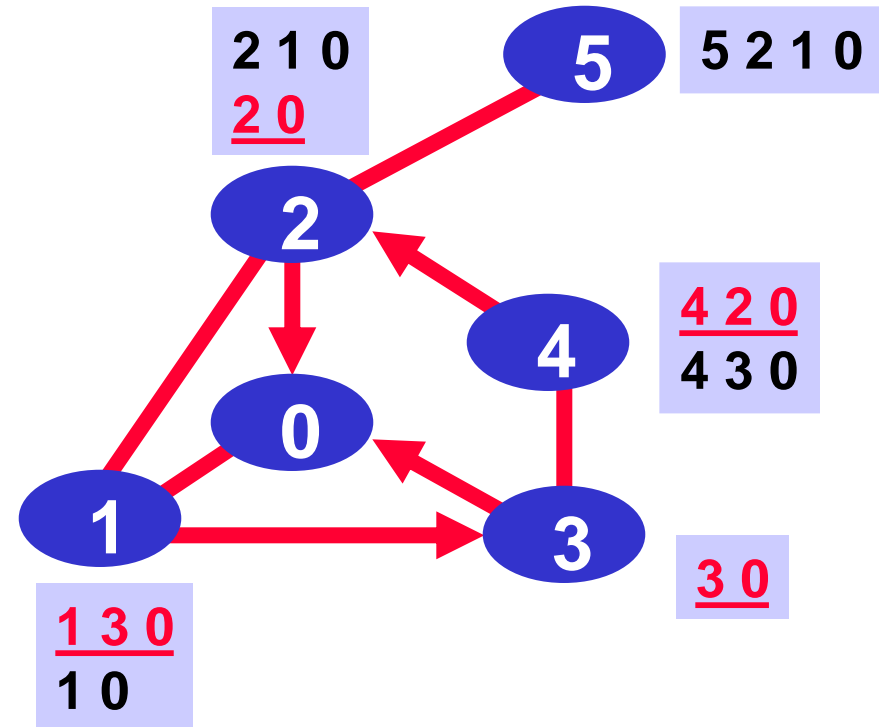
**4**

**4 2 0**
**4 3 0**

**0**

**1**

**3**

**3 0**

**1 3 0**
**1 0**

most preferred
…
least preferred (not null)

When modeling BGP : nodes represent BGP speaking routers, and 0 represents a node originating some address block

**Yes, the translation gets messy!**

# A Solution to a Stable Paths Problem

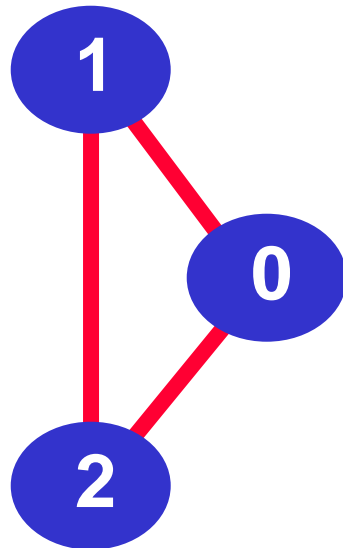A _solution_ is an assignment of permitted paths to each node such that

- node u's assigned path is either the null path or is a path uwP, where wP is assigned to node w and {u,w} is an edge in the graph,
- each node is assigned the highest ranked path among those consistent with the paths assigned to its neighbors.

**2 1 0**
**2 0**

**5 2 1 0**

**5**

**2**

**4 2 0**
**4 3 0**

**4**

**0**

**3 0**

**1**

**3**

**1 3 0**
**1 0**

A Solution need not represent a shortest path tree, or a spanning tree.
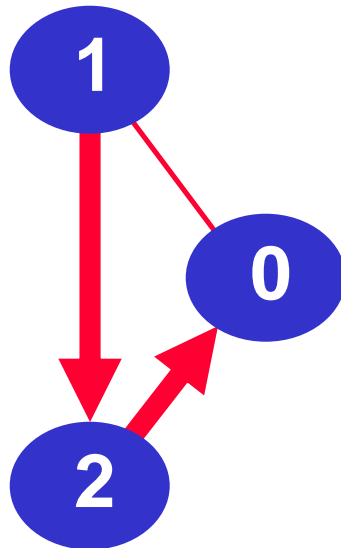
# An SPP may have multiple solutions
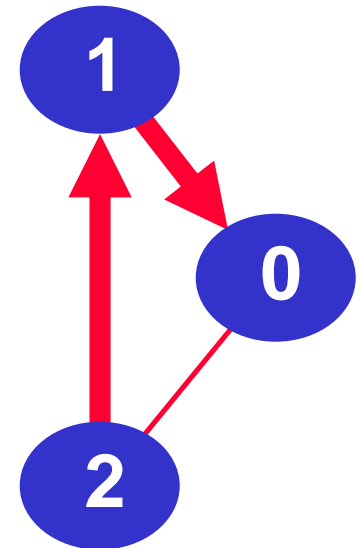
1 2 0
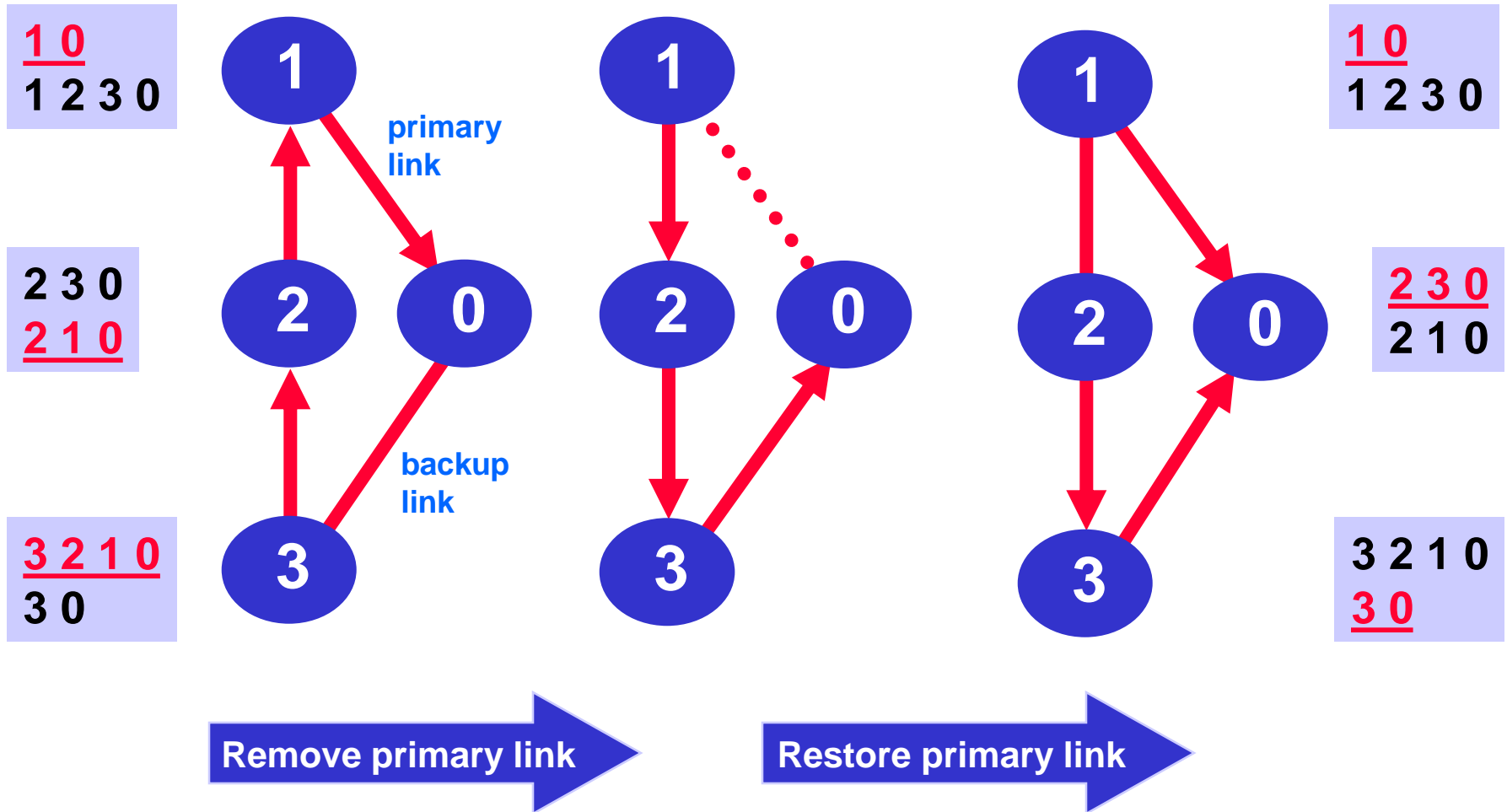1 0
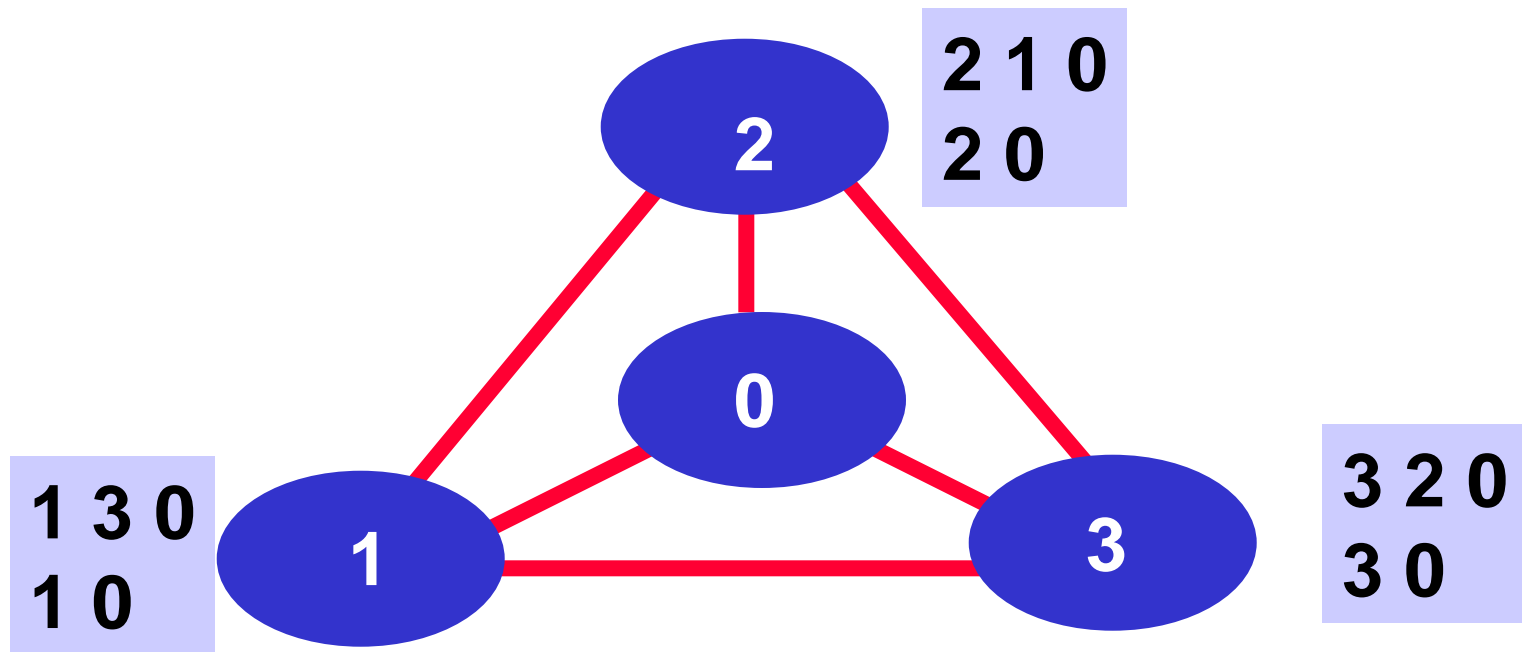
**DISAGREE**

1 2 0
1 0

2 1 0
2 0

**First solution**

1 2 0
1 0

2 1 0
2 0

**Second solution**

# Multiple solutions can result in "Route Triggering"

# BAD GADGET : No Solution



2 1 0
2 0

1 3 0
1 0

3 2 0
3 0

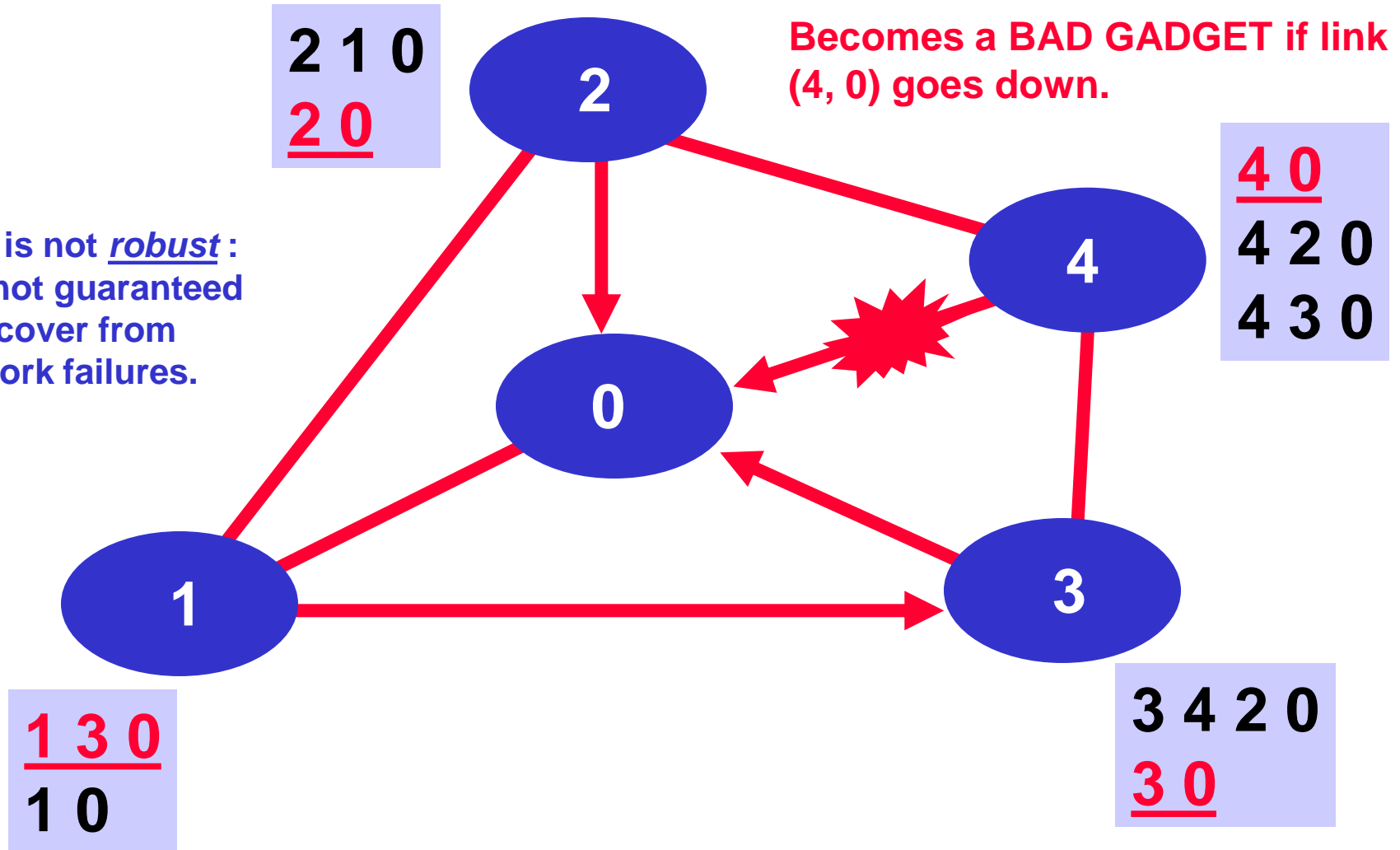**This is an SPP version of the example first presented in Persistent Route Oscillations in Inter-Domain Routing. Kannan Varadhan, Ramesh Govindan, and Deborah Estrin. Computer Networks, Jan. 2000**
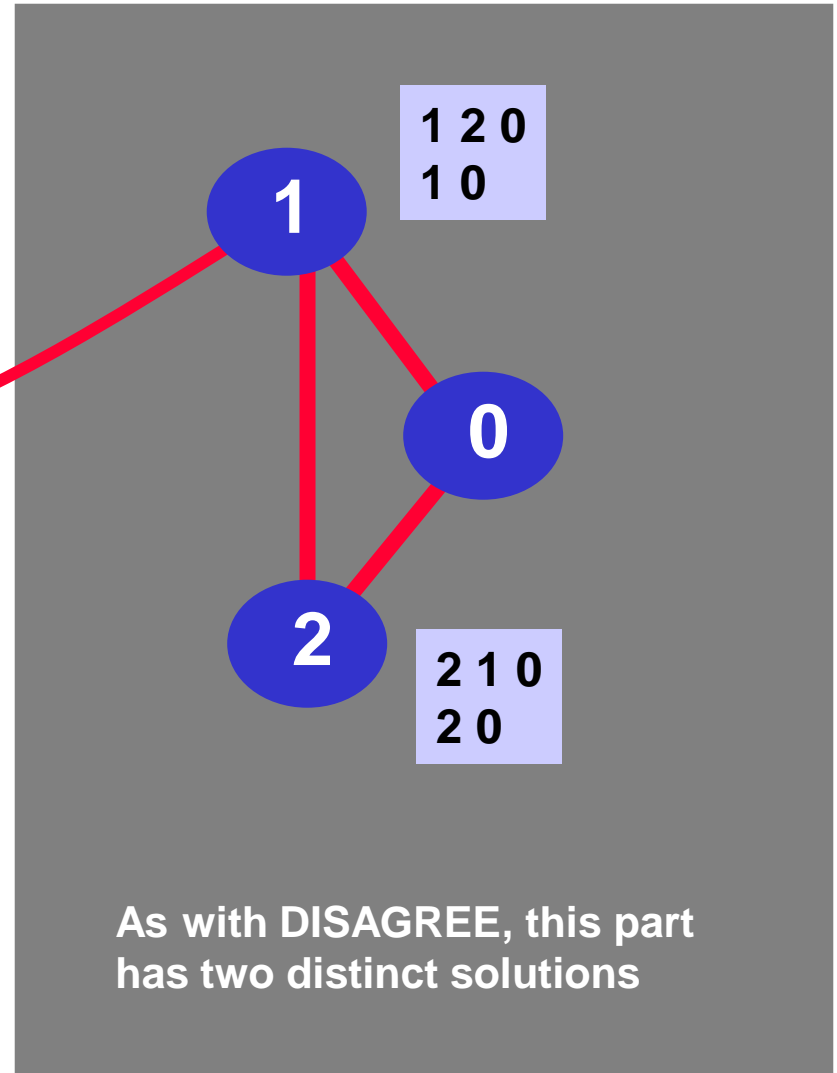
ICNP 2002

# SURPRISE : Beware of Backup Policies



**2 1 0**
**2 0**

**Becomes a BAD GADGET if link (4, 0) goes down.**

**BGP is not _robust_ :**
**it is not guaranteed**
**to recover from**
**network failures.**

**4 0**
**4 2 0**
**4 3 0**

**2**

**4**

**0**

**1**

**3**

**1 3 0**
**1 0**

**3 4 2 0**
**3 0**

# PART VI

## Current Internet Growth Trends

ICNP 2002

# Large BGP Tables Considered Harmful

- Routing tables must store best routes and alternate routes
- Burden can be large for routers with many alternate routes (route reflectors for example)
- Routers have been known to die
- Increases CPU load, especially during session reset

Moore's Law may save us in theory. But
in practice it means spending money to upgrade equipment ...

# BGP Routing Tables

```
show ip bgp
BGP table version is 111849680, local router ID is 203.62.248.4
Status codes: s suppressed, d damped, h history, * valid, > best, i - internal
Origin codes: i - IGP, e - EGP, ? - incomplete

   Network          Next Hop            Metric LocPrf Weight Path

. . .
*>i192.35.25.0      134.159.0.1              50         0 16779 1 701 703 i
*>i192.35.29.0      166.49.251.25            50         0 5727 7018 14541 i
*>i192.35.35.0      134.159.0.1              50         0 16779 1 701 1744 i
*>i192.35.37.0      134.159.0.1              50         0 16779 1 3561 i
*>i192.35.39.0      134.159.0.3              50         0 16779 1 701 80 i
*>i192.35.44.0      166.49.251.25            50         0 5727 7018 1785 i
*>i192.35.48.0      203.62.248.34            55         0 16779 209 7843 225 225 225 225 225 i
*>i192.35.49.0      203.62.248.34            55         0 16779 209 7843 225 225 225 225 225 i
*>i192.35.50.0      203.62.248.34            55         0 16779 3549 714 714 714 i
*>i192.35.51.0/25   203.62.248.34            55         0 16779 3549 14744 14744 14744 14744 14744 14744 14744 14744 i
. . .
```
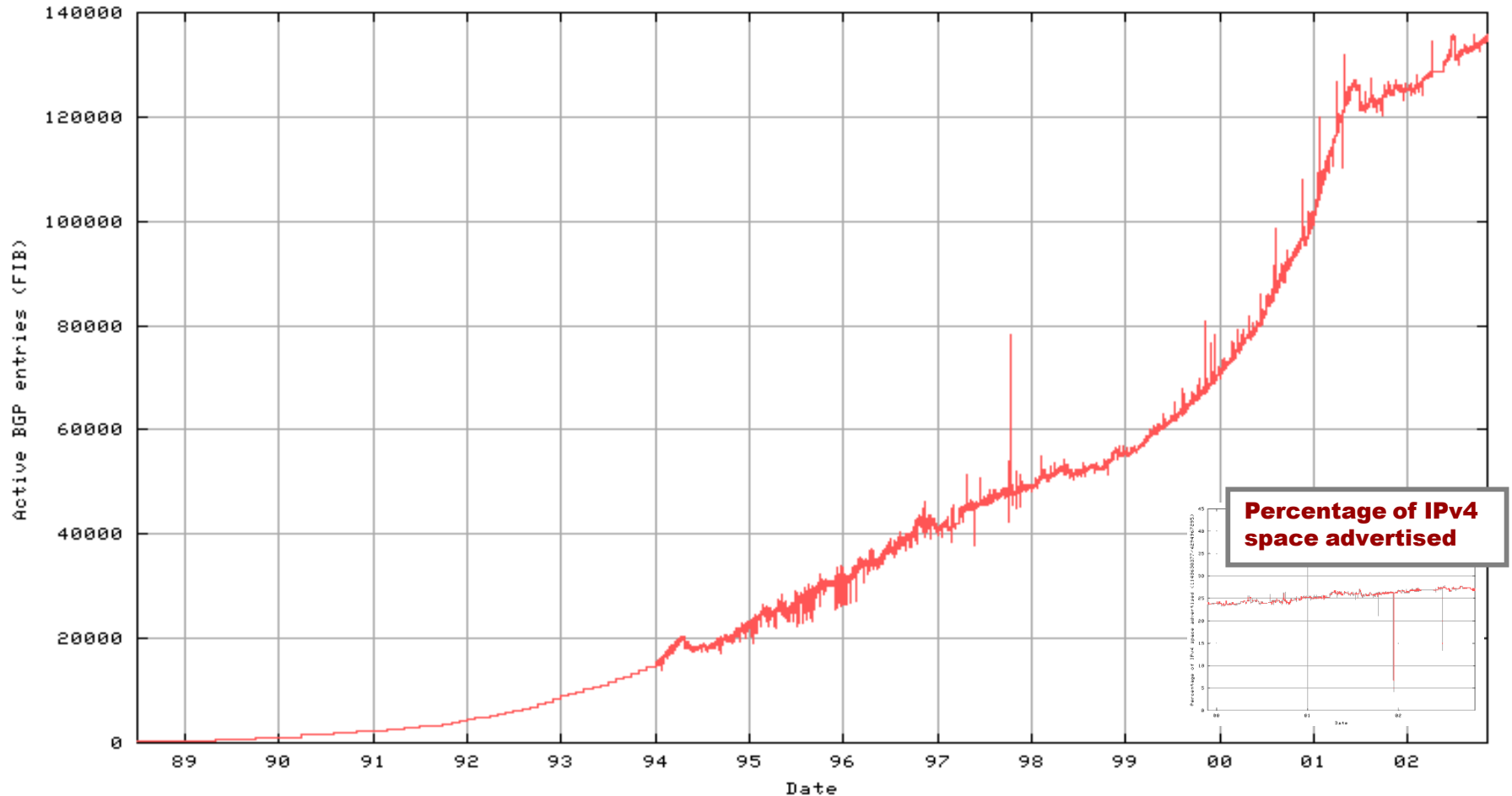
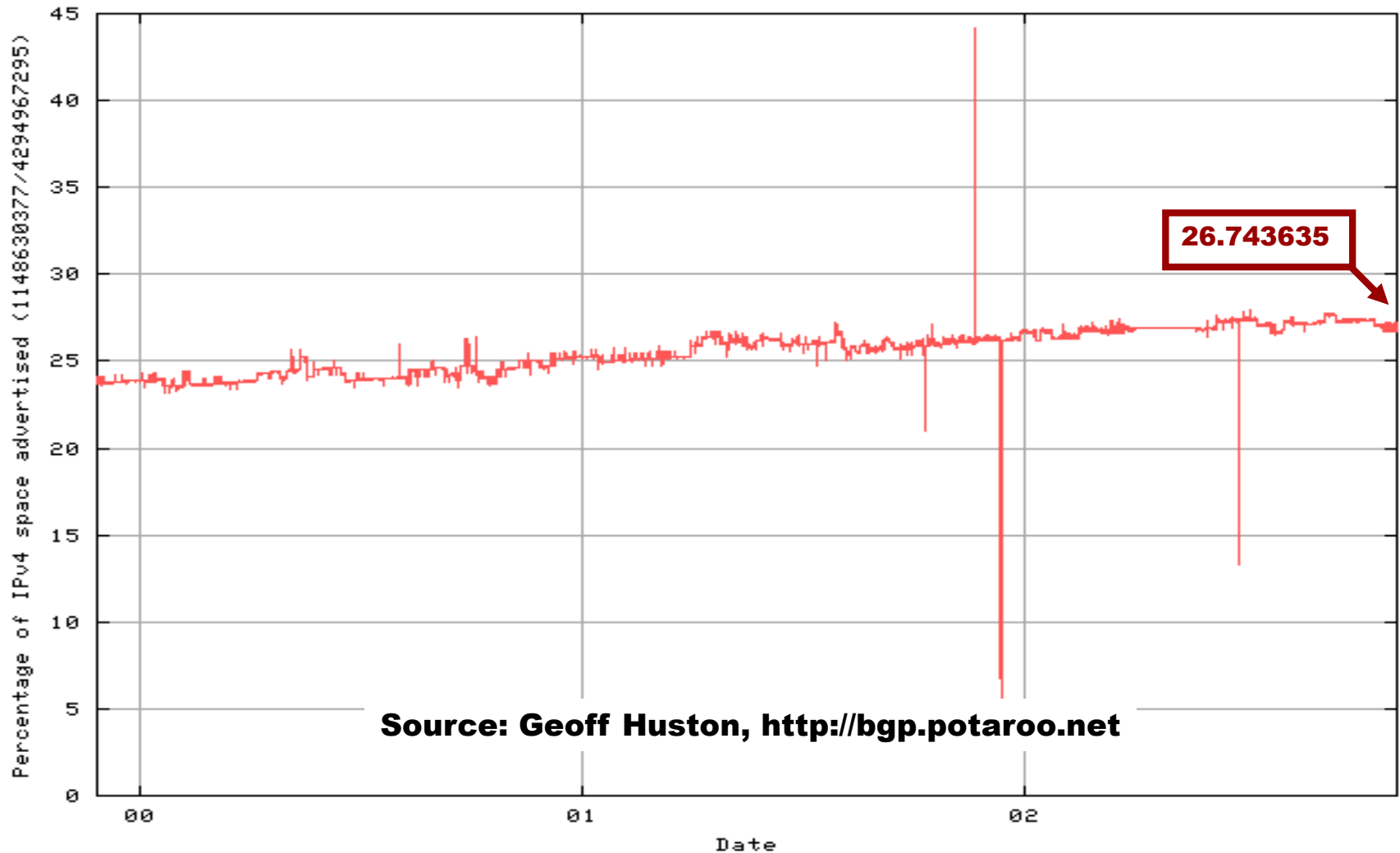**Thanks to Geoff Huston. http://www.telstra.net/ops on July 6, 2001**

- **Use "whois" queries to associate an ASN with "owner" (for example, http://www.arin.net/whois/arinwhois.html)**
- **7018 = AT&T Worldnet, 701 =Uunet,  3561 = Cable & Wireless, …**

# Growth of BGP Routes



Source: Geoff Huston, http://bgp.potaroo.net, Nov. 3, 2002

ICNP 2002

# Percent of IPv4 Space Covered



26.743635

Source: Geoff Huston, http://bgp.potaroo.net

# Average Span of BGP Prefixes



Source: Geoff Huston, http://bgp.potaroo.net

8518

# Prefix Lengths



Source: Geoff Huston, http://bgp.potaroo.net

# Number of Used ASNs



ASN Count (November 3, 2002.   Source Geoff Huston)

Source: Geoff Huston, http://bgp.potaroo.net

Legend:
- All ASNs
- Originating only
- Transit only
- Originating and Transit

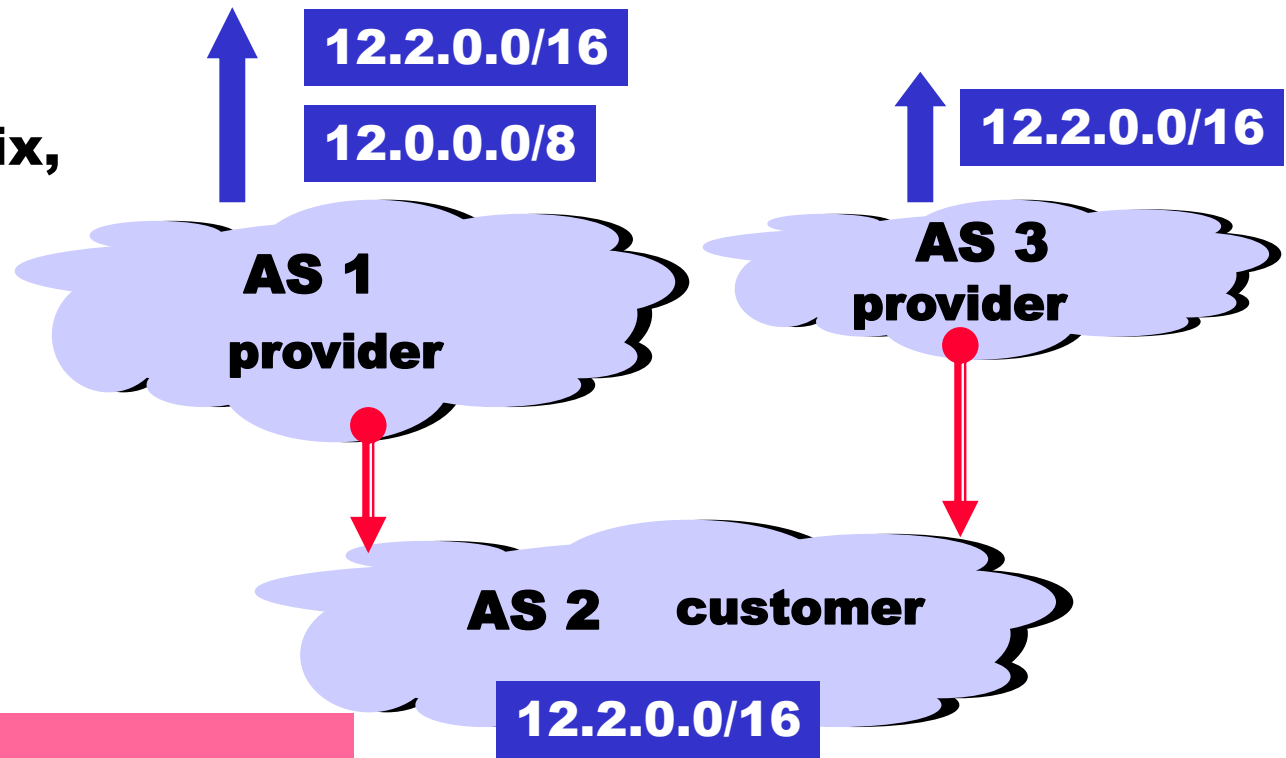# Deaggregation Due to Multihoming May Contribute to Table Growth

If AS 1 does not announce the more specific prefix, then most traffic to AS 2 will go through AS 3 because it is a longer match

12.2.0.0/16

12.0.0.0/8

12.2.0.0/16

**AS 1** provider

**AS 3** provider

**AS 2** customer

12.2.0.0/16

AS 2 is "punching a hole" in The CIDR block of AS 1

ICNP 2002

# For a Detailed Analysis ....

**Internet Expansion, Refinement, and Churn**

**Andre Broido, Evi Nemeth, and kc claffy**
**Cooperative Association for Internet Data Analysis - CAIDA**
**San Diego Supercomputer Center,**
**University of California, San Diego**

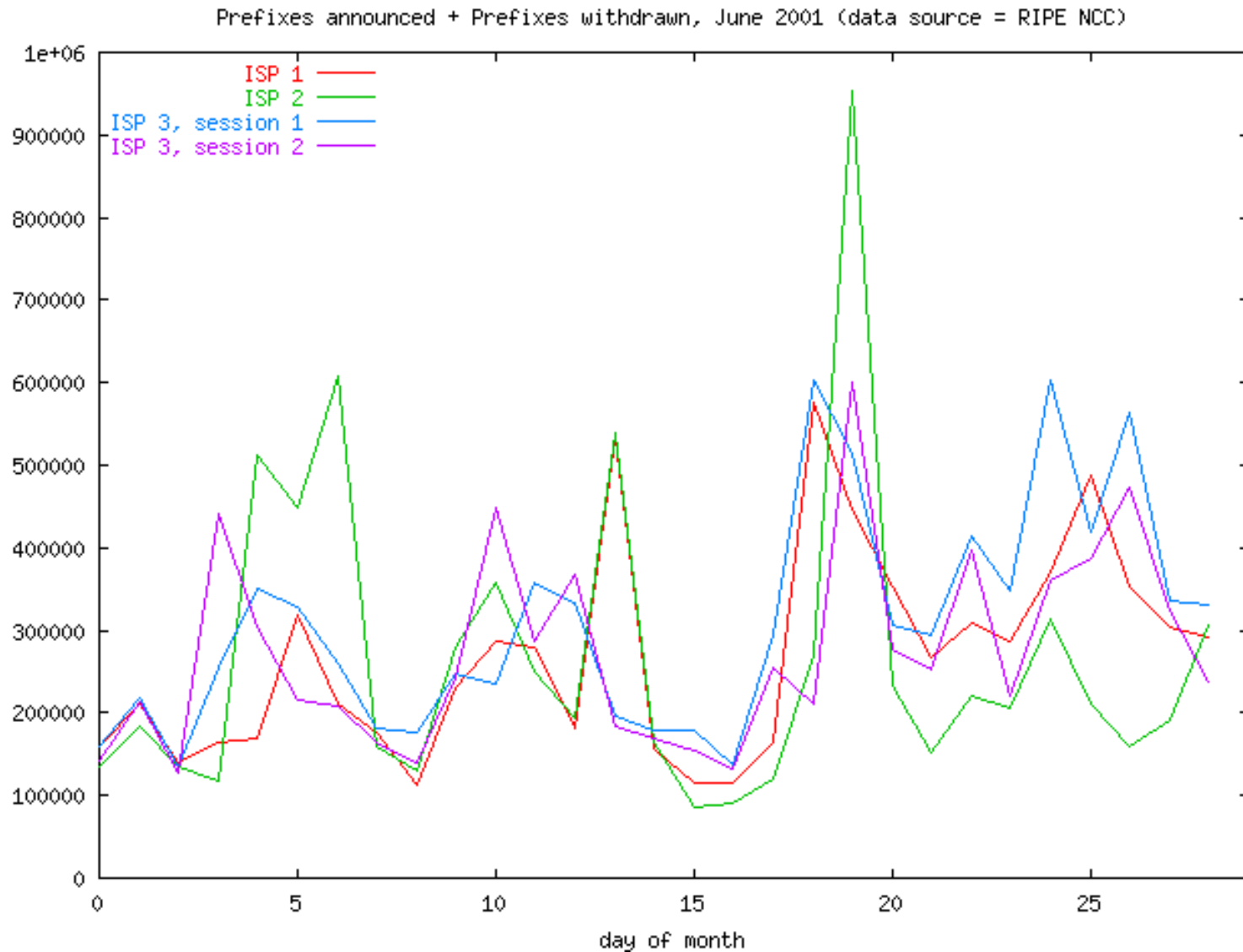**http://www.caida.org/outreach/papers/2002/EGR/**

# BGP Dynamics

- ## How many updates are flying around the Internet?

- ## How long Does it take Routes to Change?

The goals of
   (1) fast convergence
   (2) minimal updates
   (3) path redundancy
are at odds

# Daily Update Count



Prefixes announced + Prefixes withdrawn, June 2001 (data source = RIPE NCC)

# What is the Sound of One Route Flapping?



Prefixes announced + Prefixes withdrawn, June 25, 2001 (data source = RIPE NCC)

Legend: ISP 1, ISP 2, ISP 3, session 1, ISP 3, session 2

x-axis: 10 minute bins

# A Few Bad Apples ...

Cumulative Percentage on 05/28/2001

Most prefixes are
stable most of the time.
On this day, about 83% of the prefixes
were not updated.

Typically, 80% of
the updates are
for less than 5%
Of the prefixes.

ISP 1
ISP 2
ISP 3, Session 1
ISP 3, Session 2

Percent of BGP table prefixes

Data source: RIPE NCC
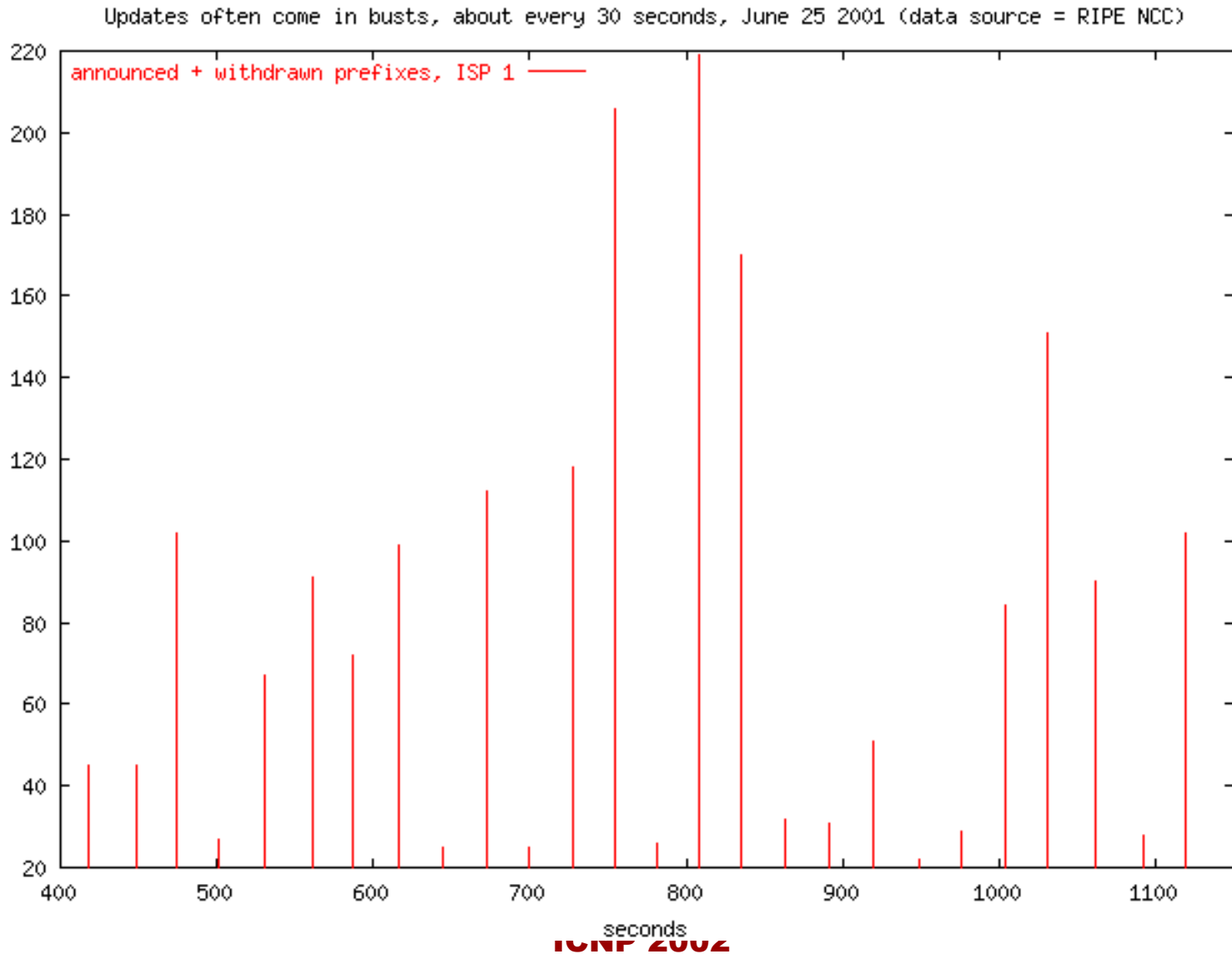
# Two BGP Mechanisms for Squashing Updates

- ## Rate limiting on sending updates
  - Send batch of updates every MinRouteAdvertisementInterval seconds (+/- random fuzz)
  - Default value is 30 seconds
  - A router can change its mind about best routes many times within this interval without telling neighbors

  **Effective in dampening oscillations inherent in the vectoring approach**

- ## Route Flap Dampening
  - Punish routes for "misbehaving"

  **Must be turned on with configuration**

# 30 Second Bursts



Updates often come in busts, about every 30 seconds, June 25 2001 (data source = RIPE NCC)

# How Long Does BGP Take to Adapt to Changes?



**Thanks to Abha Ahuja and Craig Labovitz for this plot.**

ICNP 2002

# Two Main Factors in Delayed Convergence
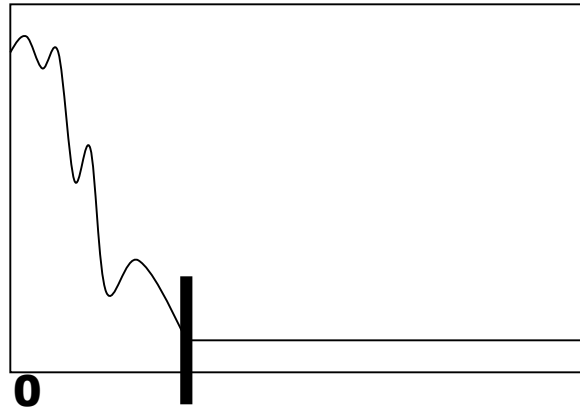
- ## Rate limiting timer slows everything down

- ## BGP can explore many alternate paths before giving up or arriving at a new path

  - ### No global knowledge in vectoring protocols

# Why is Rate Limiting Needed?

**Updates
to convergence**

**Time
to convergence**

0

0

**MinRouteAdvertisementInterval**

**MinRouteAdvertisementInterval**

## Rate limiting dampens some of the oscillation inherent in a vectoring protocol.

## Current interval (30 seconds) was picked "out of the blue sky"

**SSFNet (www.ssfnet.org) simulations, T. Griffin and B.J. Premore.
To appear in ICNP 2001.**

# Route Flap Dampening (RFC 2439)

Routes are given a <u>penalty</u> for changing. If penalty exceeds <u>suppress limit</u>, the route is dampened. When the route is not changing, its penalty decays exponentially. If the penalty goes below <u>reuse limit</u>, then it is announced again.

- Can dramatically reduce the number of BGP updates
- Requires additional router resources
- Applied on eBGP inbound only

# Route Flap Dampening Example



Route Flap Dampening (halflife = 15 minutes) -- 6 flaps, one every 10 minutes

penalty for each flap = 1000

ICNP 2002

# Q: Why All the Updates?

- Networks come, networks go
- There's always a router rebooting somewhere
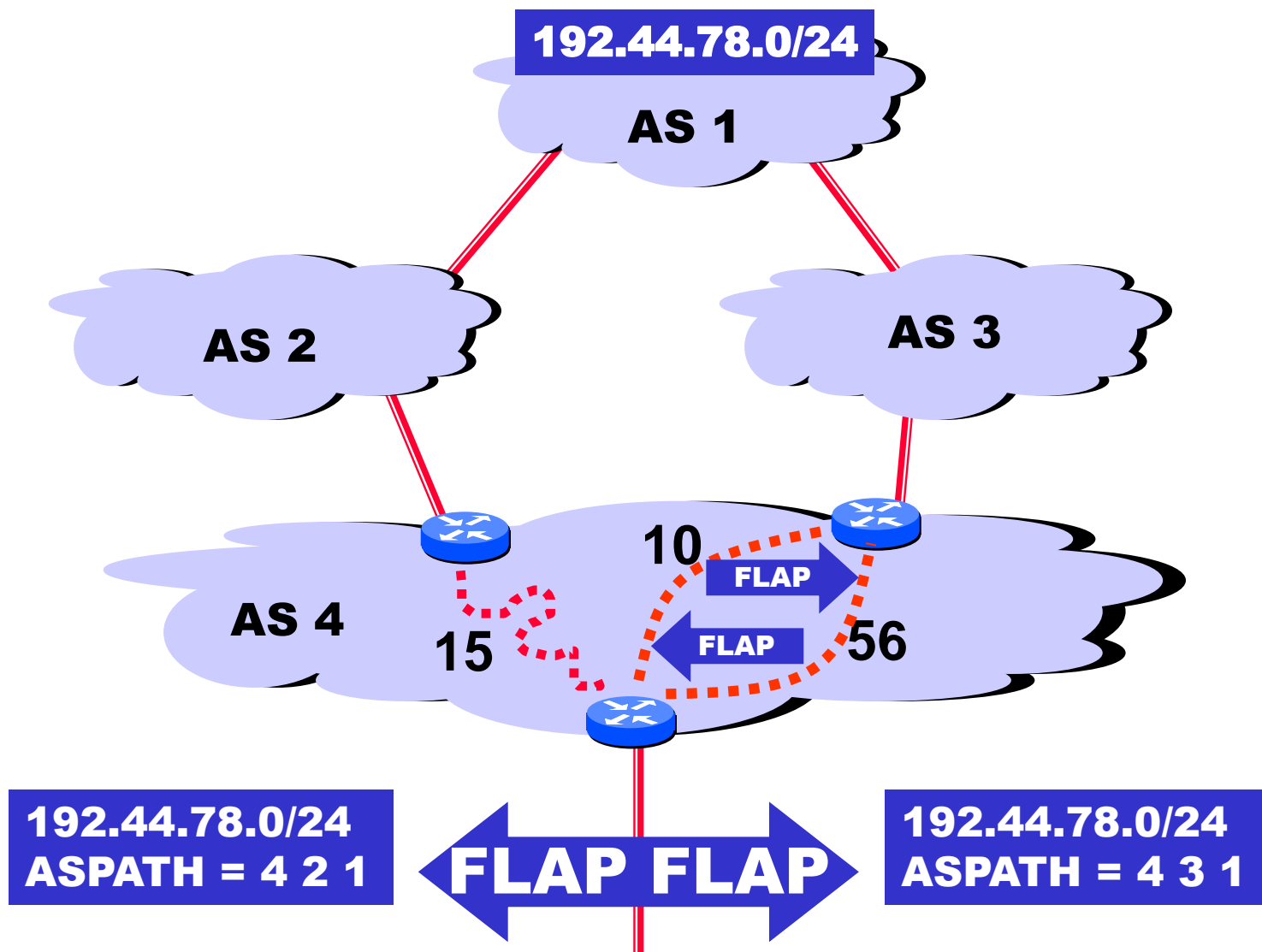- Hardware failure, flaky interface cards, backhoes digging, floods in Houston, ...

This is "normal" --- exactly what dynamic routing is designed for...

# Q: Why All the Updates?

- Misconfiguration
- Route flap dampening not widely used
- BGP exploring many alternate paths
- Software bugs in implementation of routing protocols
- BGP session resets due to congestion or lack of interoperability: BGP sessions are brittle. One malformed update is enough to reset session and flap 100K routes. (Consequence of incremental approach)
- IGP instability exported by use of MEDs or IGP tie breaker
- Sub-optimal vendor implementation choices
- Secret sauce routing algorithms attempting fancy-dancy tricks
- Weird policy interactions (MED oscillation, BAD GADGETS??)
- Gnomes, sprites, and fairies
- ....

**A: NO ONE <u>REALLY</u> KNOWS ...**

# IGP Tie Breaking Can Export Internal Instability to the Whole Wide World

192.44.78.0/24

AS 1

AS 2

AS 3

AS 4

10

FLAP

FLAP

15

56

192.44.78.0/24
ASPATH = 4 2 1

**FLAP FLAP**

192.44.78.0/24
ASPATH = 4 3 1

# MEDs Can Export Internal Instability



114

# Implementation Does Matter!



**Thanks to Abha Ahuja and Craig Labovitz for this plot.**
ICNP 2002

# How Long Will Interdomain Routing Continue to Scale?

**A quote from some recent email:**

> ... the existing interdomain routing infrastructure is rapidly nearing the end of its useful lifetime. It appears unlikely that mere tweaks of BGP will stave off fundamental scaling issues, brought on by growth, multihoming and other causes.

**Is this true or false?    How can we tell?**

**Research required...**

# Summary

- BGP is a fairly simple protocol …
- … but it is not easy to configure
- BGP is running on more than 100K routers (my estimate), making it one of world's largest and most visible distributed systems
- Global dynamics and scaling principles are still not well understood

# PART VII

## Selected Bibliography

# Addressing and ASN RFCs

- RFC 1380 IESG Deliberations on Routing and Addressing (1992)
-   RFC 1517Applicability Statement for the Implementation of Classless Inter-Domain Routing (CIDR) (1993)
-   RFC 1518 An Architecture for IP Address Allocation with CIDR (1993)
-   RFC 1519 Classless Inter-Domain Routing (CIDR) (1993)
-   RFC 1467 Status of CIDR Deployment in the Intrenet (1983)
-   RFC 1520 Exchanging Routing Information Across Provider Boundaries in the CIDR Environment (1993)
-   RFC 1817 CIDR and Classful routing (1995)
-   RFC 1918 Address Allocation for Private Internets (1996)
-   RFC 2008 Implications of Various Address Allocation Policies for Internet Routing (1996)
-   RFC 2050 Internet Registry IP Allocation Guidelines (1996)
-   RFC 2260 Scalable Support for Multi-homed Multi-provider Connectivity (1998)
-   RFC 2519 A Framework for Inter-Domain Route Aggregation (1999)
- RFC 1930 Guidelines for creation, selection, and registration of an Autonomous System (AS)
- RFC 2270 Using a Dedicated AS for Sites Homed to a Single Provider

# Selected BGP RFCs

**Internet Engineering Task Force (IETF)**   http://www.ietf.org

- IDR : http://www.ietf.org/html.charters/idr-charter.html
- RFC 1771 A Border Gateway Protocol 4 (BGP-4)
  - Latest draft rewrite: draft-ietf-idr-bgp4-18.txt
- RFC 1772 Application of the Border Gateway Protocol in the Internet
- RFC 1773 Experience with the BGP-4 protocol
- RFC 1774 BGP-4 Protocol Analysis
- RFC 2796 BGP Route Reflection An alternative to full mesh IBGP
- RFC 3065 Autonomous System Confederations for BGP
- RFC 1997 BGP Communities Attribute
- RFC 1998 An Application of the BGP Community Attribute in Multi-home Routing
- RFC 2439 Route Flap Dampening

# Titles of Some Recent Internet Drafts

- Dynamic Capability for BGP-4
- Application of Multiprotocol BGP-4 to IPv4 Multicast Routing
- Graceful Restart mechanism for BGP
- Cooperative Route Filtering Capability for BGP-4
- Address Prefix Based Outbound Route Filter for BGP-4
- Aspath Based Outbound Route Filter for BGP-4
- Architectural Requirements for Inter-Domain Routing in the Internet
- BGP support for four-octet AS number space
- Autonomous System Number Substitution on Egress
- BGP Extended Communities Attribute
- Controlling the redistribution of BGP routes
- BGP Persistent Route Oscillation Condition
- Benchmarking Methodology for Basic BGP Convergence
- Terminology for Benchmarking External Routing Convergence Measurements

## BGP is a moving target ...

ICNP 2002

# Selected Bibliography on Routing

- Internet Routing Architectures. Bassam Halabi. Second edition Cisco Press, 2000

- BGP4: Inter-domain Routing in the Internet. John W. Stewart, III.  Addison-Wesley, 1999

- Routing in the Internet. Christian Huitema. 2000

- ISP Survival Guide: Strategies for Running a Competitive ISP. Geoff Huston. Wiley, 1999.

- Interconnection, Peering and Settlements. Geoff Huston. The Internet Protocol Journal.  March and June 1999.

# BGP Stability and Convergence

- Route Flap Damping Exacerbates Internet Routing Convergence. Z.M.Mao, R.Govindan, G.Varghese,R.H.Kranz. SIGCOMM 2002.
- The Impact of Internet Policy and Topology on Delayed Routing Convergence. Craig Labovitz, Abha Ahuja, Roger Wattenhofer, Srinivasan Venkatachary. INFOCOM 2001
- An Experimental Study of BGP Convergence. Craig Labovitz, Abha Ahuja, Abhijit Abose, Farnam Jahanian. SIGCOMM 2000
- Origins of Internet Routing Instability. C. Labovitz, R. Malan, F. Jahanian. INFOCOM 1999
- Internet Routing Instability. Craig Labovitz, G. Robert Malan and Farnam Jahanian. SIGCOMM 1997

# Analysis of Interdomain Routing

- **Cooperative Association for Internet Data Analysis (CAIDA)**
  - **http://www.caida.org/**
  - **Tools and analyses promoting the engineering and maintenance of a robust, scalable global Internet infrastructure**
- **Internet Performance Measurement and Analysis (IPMA)**
  - **http://www.merit.edu/ipma/**
  - **Studies the performance of networks and networking protocols in local and wide-area networks**
- **National Laboratory for Applied Network Research (NLANR)**
  - **http://www.nlanr.net/**
  - **Analysis, tools, visualization.**
- **IRTF Routing Research Group (IRTF-RR)**
  - **http://puck.nether.net/irtf-rr/**
- **Geoff Huston: http://bgp.potaroo.net**

# Internet Route Registries

- **Internet Route Registry**
  - http://www.irr.net/
- **Routing Policy Specification Language (RPSL)**
  - RFC 2622 Routing Policy Specification Language (RPSL)
  - RFC 2650 Using RPSL in Practice
- **Internet Route Registry Daemon (IRRd)**
  - http://www.irrd.net/
- **RAToolSet**
  - http://www.isi.edu/ra/RAToolSet/

ICNP 2002

# Some BGP Theory

- **Persistent Route Oscillations in Inter-Domain Routing. Kannan Varadhan, Ramesh Govindan, and Deborah Estrin. Computer Networks, Jan. 2000. (Also USC Tech Report, Feb. 1996)**
  - Shows that BGP is not guaranteed to converge
- **An Architecture for Stable, Analyzable Internet Routing. Ramesh Govindan, Cengiz Alaettinoglu, George Eddy, David Kessens, Satish Kumar, and WeeSan Lee. IEEE Network Magazine, Jan-Feb 1999.**
  - Use RPSL to specify policies. Store them in registries.  Use registry for conguration generation and analysis.
- **An Analysis of BGP Convergence Properties. Timothy G. Griffin,  Gordon Wilfong. SIGCOMM 1999**
  - Model BGP, shows static analysis of divergence in policies is NP complete
- **Policy Disputes in Path Vector Protocols. Timothy G. Griffin,  F. Bruce Shepherd,  Gordon Wilfong. ICNP 1999**
  - Define Stable Paths Problem and develop sufficient condition for "sanity"
- **A Safe Path Vector Protocol. Timothy G. Griffin,  Gordon Wilfong. INFOCOM 2001**
  - Dynamic solution for SPVP based on histories
- **Stable Internet Routing without Global Coordination. Lixin Gao,  Jennifer Rexford. SIGMETRICS 2000**
  - Show that if certain guidelines are followed, then all is well.
- **Inherently safe backup routing with BGP. Lixin Gao, Timothy G. Griffin,  Jennifer Rexford. INFOCOM 2001**
  - Use SPP to study complex backup policies
- **On the Correctness of IBGP Configurations. Griffin and Wilfong.SIGCOMM 2002.**
- **An Analysis of the MED oscillation Problem. Griffin and Wilfong.  ICNP 2002.**